

APPENDIX

TEST: TEXT PROTOTYPE ALIGNED EMBEDDING TO ACTIVATE LLM’S ABILITY FOR TIME SERIES

Anonymous authors

Paper under double-blind review

1 RELATED WORK

Our work mainly involves two research fields: Universal Representation Learning (URL) for time series based on Contrastive Learning (CL) and Large Language Model (LLM) + Time Series (TS).

1.1 CL-BASED URL FOR TS

Unsupervised URL approaches aim to learn discriminative feature representations from unlabeled data, without the requirement of annotating every sample. Enabling URL is extremely crucial for time series data, due to its unique annotation bottleneck caused by its complex characteristics and lack of visual cues compared with other data modalities.

Contrastive methods learn meaningful representations from time series by optimizing self-discrimination tasks. Instead of directly modeling the complex raw data, they employ pretext tasks that leverage the underlying similarity between samples, which eliminates the need for re-constructing the complete input and allows for the discovery of contextualized underlying factors of variations. Contrastive methods typically generate augmented views of the raw data through various transformations and then learn representations by contrasting positive samples against negative samples. The existing CL-based URL for TS are listed in Table S1.

Type	Methods			
Instance-level	SimCLR Chen et al. (2020)	TimeCLR Yang et al. (2022)	MoCo He et al. (2020)	BYOL Grill et al. (2020)
	CPC van den Oord et al. (2018)	SimSiam Zheng et al. (2023)	MCL Wickstrøm et al. (2022)	
Prototype-level	SwAV Caron et al. (2020)	PCL Li et al. (2021b)	CCL Sharma et al. (2020)	SCCL Zhang et al. (2021)
	CC Li et al. (2021c)	SLIC Khorasgani et al. (2022)	MHCCL Meng et al. (2022)	
Temporal-level	TS2Vec Yue et al. (2022)	TS-TCC Eldele et al. (2021)	TNC Tonekaboni et al. (2021)	TCL
	T-Loss Franceschi et al. (2019b)	BTSF Yang & Hong (2022)	CoST Woo et al. (2022a)	

Table S1: Contrastive Learning based Universal Representation Methods for Time Series

Instance-level contrastive models treat individual samples independently for the purpose of instance discrimination. They utilize data augmentations to transform original inputs into a new embedding space. Within this space, augmentations derived from the same sample are considered as positive pairs, while those from different samples are treated as negative pairs. During training, these models are optimized by maximizing the similarity between representations of positive pairs, while simultaneously minimizing the similarity between representations of negative pairs.

Prototype-level contrastive models break the independence between samples and explore to exploit the implicit semantics shared by samples in the same cluster. They can address the limitation that instance-level contrastive learning models tend to treat semantically similar samples as negatives.

Temporal-level contrastive models instead focus on capturing scale- invariant representations at each individual timestamp. By considering both instance-level and temporal-level representation learning

strategies, researchers aim to enhance the capability of contrastive learning methods in capturing the complexities inherent in time series data.

1.2 LLM+TS

Large models, specifically referred to as large language models (LLMs) and pre-trained foundation models (PFMs), have witnessed remarkable success across a multitude of tasks and domains, such as natural language processing (NLP), computer vision (CV). Given the remarkable achievements of large models in these diverse fields, an intriguing question emerges: can large models be effectively employed to analyze TS data?

TS data has long been studied and proven to be indispensable in a myriad of real-world applications, encompassing fields such as geoscience, transportation, energy, healthcare, environment, and finance. While large models have made significant progress in various fields, the arena of time series analysis has followed a more gradual path. Traditional analytical methods have predominantly relied on statistical models. The advent of deep learning has galvanized the research community to explore more potent data-driven models, typically built on the basis of Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and Transformers. Nonetheless, the majority of these models remain relatively small in scale and are tailored for specific tasks, thereby lacking the capacity to acquire comprehensive semantic and knowledge representations from large-scale data for multi-task reasoning.

There hasn't been much research done on TS+LLM because this field is still in its infancy. We summarize the existing work in Table S2. Different from the main text, we category work here through technical means.

Means	Pros	Cons	Work
Training	Specialized, accurate	Not universal, large datasets	Pre-training Ma et al. (2023) Earth transformer Bi et al. (2023)
Tuning	End-to-end, accurate	More experiments, lose language ability	GPT4TSZhou et al. (2023) LLM4TSChang et al. (2023) LLMTime Gruver et al. (2023)
Tool Augmented	Parameter-efficient, less experiments	Need experts, need annotation	PromptCast Xue & Salim (2023) Health Learner Liu et al. (2023) METS Li et al. (2023) Text2ECGChung et al. (2023)
External Encoder	Parameter-efficient, multiple abilities	Weak robust	TEST Time-LLM Jin et al. (2023)

Table S2: Existing Work about TS+LLM

2 EXPERIMENTS

2.1 MODEL

2.1.1 ENCODER

The core of TEST is to train an encoder and a soft prompt. The encoder must can extract relevant information from TS, needs to be time- and memory-efficient, and has to allow variable-length inputs. Thus, as shown in Figure S1, we build a causal TCN with 10 layers of convolution blocks. Each convolution block is a sequence of GELU, DilatedConv, BatchNorm, GELU, DilatedConv, with skip connections across each block. The DilatedConvs have dilation of $2i$ in each layer i of convolution block. A final convolution block is used to map the hidden channels to the output channel whose size is the same as the LLM's embedding size.

The detailed architecture is: Number of channels in the intermediary layers of the causal network is 40; Number of layers (depth of the causal network) is 10; Kernel size of all convolutions is 3; Negative slope of the leaky ReLU activation is 0.01; Number of output channels of the causal network (before max pooling) is 640; Dimension of the representations is the same as the LLM's embedding size (e.g. 1024 for gpt2).

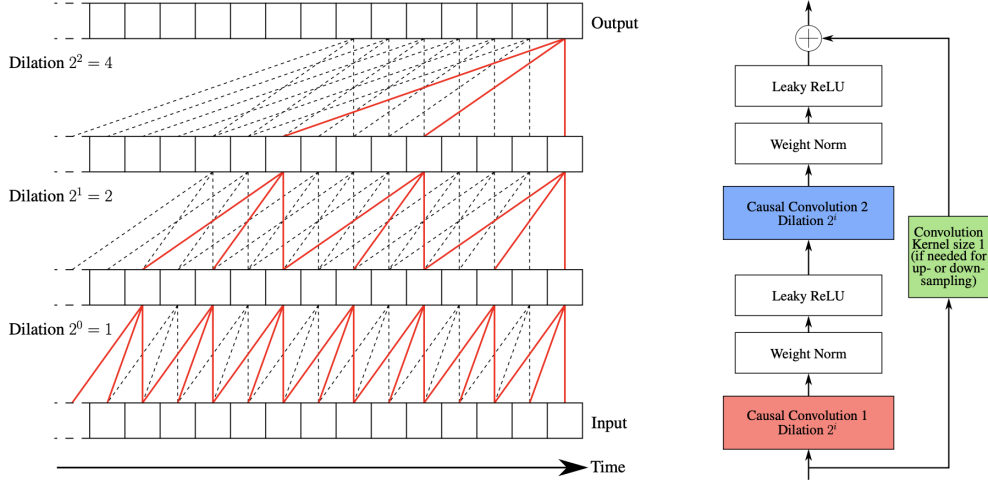


Figure S1: Illustration of Three Stacked Dilated Causal Convolutions and Composition of the i -th Layer of The Chosen Architecture

We train our models with the following parameters for time series classification. Note that no hyperparameter optimization was performed on the encoder hyperparameters: Optimizer is Adam with learning rate $\alpha = 0.001$ and decay rates $\beta = (0.9, 0.999)$; Number of negative samples is $K \in \{1, 2, 5, 10\}$ for univariate time series, $K \in \{5, 10, 20\}$ for multivariate ones; Batch size is 10; Number of optimization steps is 2000 for $K \leq 10$ (i.e., 20 epochs for a dataset of size 1000), 1500 otherwise.

2.1.2 LLM

The used LLMs are as listed in Table S3. Each encoder and soft prompt of LLM are trained using the Adam optimizer on 20 NVIDIA Tesla V100-SXM2 GPU with CUDA 11.3.

Model	Size	Embed. dimension
Bert Devlin et al. (2018)	110M, 335M	748, 1024
GPT2 Radford et al. (2019)	117M, 345M, 774M	768, 1024, 1280
ChatGLM Du et al. (2022)	6B	4096
LLaMa2 Touvron et al. (2023)	7B, 13B	4096

Table S3: The Used Language Model

2.2 FORECASTING TASKS

All the deep learning networks are implemented in PyTorch and trained on NVIDIA V100 32GB GPUs. We use mean square error (MSE) and mean absolute error (MAE) as metrics. For zero-shot learning, mean absolute percentage error (MAPE) is used for TOURISM; symmetric MAPE (sMAPE) is used for M3 and M4; normalized deviation (ND) is used for ELECTR. All experiments are repeated 3 times and the mean of the metrics is used in the final results.

2.2.1 DATASET DETAILS

The details of long-term forecasting and few-shot forecasting datasets are: ETT datasets Zhou et al. (2021) contain electricity load of various resolutions (ETTh & ETTm) from two electricity stations; Weather dataset Wetterstation (2017) contains 21 meteorological indicators of Germany within 1 year; Illness dataset CDC (2021) contains the influenza-like illness patients in the United States. ILI is not used for few-shot learning for the limited quantity that is hard to follow the definition of few-shot; Electricity dataset SJ & B (2017) contains the electricity consumption; Traffic dataset

PeMS (2021) contains the occupation rate of freeway system across the State of California. Table S4 summarizes details of feature statistics.

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

Table S4: Long-term Forecasting and Few-shot Forecasting Dataset Details

	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	18	-	Yearly
M4 Quarterly	6	24000	-	Quarterly
M4 Monthly	8	48000	-	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

Table S5: Zero-term Forecasting Datasets and Mapping Details of Zero-shot Learning

The details of zero-shot forecasting datasets are: M4 is a large and diverse dataset that contains time series of various frequencies and fields, including business, financial and economic forecasting; M3 is smaller than M4, but also contains time series from diverse domains and frequencies; TOURISM is the dataset of tourism activities with different frequencies and contains a much higher fraction of erratic series compared with M4; ELECTR represents the electricity usage monitoring of 370 customers over three years. Table S5 summarizes details of the datasets and zero-shot mapping between source and target.

2.2.2 BASELINE DETAILS

For long-shot forecasting, we refer to the SOTA methods reported in Wu et al. (2023): TimesNet Wu et al. (2023), ETSformer Woo et al. (2022b), DLinear Zeng et al. (2023), FEDformer Zhou et al. (2022), Informer Zhou et al. (2021), and LLM for TS method GPT4TS Zhou et al. (2023).

For few-shot forecasting, we refer to the SOTA methods reported in Zhou et al. (2023): DLinear Zeng et al. (2023), PatchTST Nie et al. (2023), TimesNet Wu et al. (2023), FEDformer Zhou et al. (2022), Autoformer Wu et al. (2021), Stationary Liu et al. (2022), ETSformer Woo et al. (2022b), Informer Zhou et al. (2021), Reformer Kitaev et al. (2020)

For zero-shot forecasting, we refer to the SOTA methods reported in Zhou et al. (2023): N-BEATS Oreshkin et al. (2020), DLinear Zeng et al. (2023), PatchTST Nie et al. (2023), TimesNet Wu et al. (2023), FEDformer Zhou et al. (2022), Autoformer Wu et al. (2021), Stationary Liu et al. (2022), ETSformer Woo et al. (2022b), Informer Zhou et al. (2021), Reformer Kitaev et al. (2020)

2.2.3 LONG-TERM FORECASTING

We follow the classical experiment settings and the results of SOTA models in Wu et al. (2023) (ICLR 2023). The results are shown in Table S6. Overall, TEST achieves comparable performance to SOTA models TimesNet and Dlinear, and outperforms other baselines.

Methods		TEST	GPT4TS	TimesNet	ETSformer	DLinear	FEDformer	Informer	TCN	LSTM
ETTh1	96	0.293 0.346	0.292 0.346	0.325 0.398	0.338 0.375	0.345 0.372	0.375 0.398	0.672 0.571	0.863 0.664	0.863 0.664
	192	0.332 0.369	0.332 0.372	0.324 0.387	0.408 0.410	0.380 0.389	0.426 0.441	0.795 0.669	0.837 0.700	1.113 0.776
	336	0.368 0.392	0.366 0.394	0.360 0.411	0.435 0.428	0.413 0.413	0.445 0.459	1.212 0.871	1.124 0.832	1.267 0.832
	720	0.418 0.420	0.417 0.421	0.428 0.450	0.499 0.462	0.474 0.453	0.543 0.490	1.166 0.823	1.153 0.820	1.324 0.858
	Avg	0.353 0.382	0.352 0.383	0.350 0.406	0.429 0.425	0.403 0.407	0.448 0.452	0.961 0.734	0.929 0.725	1.142 0.782
ETTh2	96	0.372 0.400	0.376 0.397	0.384 0.402	0.494 0.479	0.386 0.400	0.376 0.419	0.865 0.713	0.878 0.740	1.044 0.773
	192	0.414 0.422	0.416 0.418	0.436 0.429	0.538 0.504	0.437 0.432	0.420 0.448	1.008 0.792	1.037 0.824	1.217 0.832
	336	0.422 0.437	0.442 0.433	0.491 0.469	0.574 0.521	0.481 0.459	0.459 0.465	1.107 0.809	1.238 0.932	1.259 0.841
	720	0.447 0.467	0.477 0.456	0.521 0.500	0.562 0.535	0.519 0.516	0.506 0.507	1.181 0.865	1.135 0.852	1.271 0.838
	Avg	0.414 0.431	0.427 0.426	0.458 0.450	0.542 0.510	0.456 0.452	0.440 0.460	1.040 0.795	1.072 0.837	1.198 0.821
Electricity	96	0.275 0.338	0.285 0.342	0.340 0.374	0.340 0.391	0.333 0.387	0.358 0.397	3.755 1.525	2.116 1.197	2.522 1.278
	192	0.340 0.379	0.354 0.389	0.402 0.414	0.430 0.439	0.477 0.476	0.429 0.439	5.602 1.931	4.315 1.635	3.312 1.384
	336	0.329 0.381	0.373 0.407	0.452 0.452	0.485 0.559	0.594 0.541	0.496 0.487	4.721 1.835	1.124 1.604	3.291 1.388
	720	0.381 0.423	0.406 0.441	0.462 0.468	0.500 0.497	0.831 0.657	0.463 0.474	3.647 1.625	3.188 1.540	3.257 1.357
	Avg	0.331 0.380	0.354 0.394	0.414 0.427	0.439 0.452	0.559 0.515	0.4370.449	4.431 1.729	2.686 1.494	3.095 1.352
Traffic	96	0.132 0.223	0.139 0.238	0.168 0.222	0.187 0.304	0.197 0.282	0.193 0.308	0.274 0.368	0.258 0.357	0.375 0.437
	192	0.158 0.241	0.153 0.251	0.184 0.239	0.199 0.196	0.285 0.201	0.315 0.296	0.386 0.266	0.368 0.348	0.442 0.473
	336	0.163 0.260	0.169 0.266	0.198 0.260	0.212 0.329	0.209 0.301	0.214 0.329	0.300 0.394	0.280 0.380	0.439 0.473
	720	0.199 0.291	0.206 0.297	0.220 0.300	0.233 0.345	0.245 0.333	0.246 0.355	0.373 0.439	0.283 0.376	0.980 0.814
	Avg	0.162 0.253	0.167 0.263	0.192 0.245	0.208 0.323	0.212 0.300	0.214 0.327	0.311 0.397	0.313 0.401	0.559 0.549
Weather	96	0.407 0.282 0	0.388 0.282	0.593 0.321	0.607 0.392	0.650 0.396	0.587 0.366	0.719 0.391	0.684 0.384	0.843 0.453
	192	0.423 0.287	0.407 0.290	0.617 0.336	0.621 0.399	0.598 0.370	0.604 0.373	0.696 0.379	0.685 0.390	0.847 0.453
	336	0.430 0.296	0.412 0.294	0.629 0.336	0.622 0.396	0.605 0.373	0.621 0.383	0.777 0.420	0.734 0.408	0.853 0.455
	720	0.463 0.315	0.450 0.312	0.640 0.350	0.632 0.396	0.645 0.394	0.626 0.382	0.864 0.472	0.717 0.396	1.500 0.805
	Avg	0.430 0.295	0.414 0.294	0.620 0.336	0.621 0.396	0.625 0.383	0.610 0.376	0.764 0.416	0.705 0.395	1.011 0.541
ILI	96	0.150 0.202	0.162 0.212	0.152 0.220	0.197 0.281	0.196 0.255	0.217 0.296	0.300 0.384	0.458 0.490	0.369 0.406
	192	0.198 0.246	0.204 0.248	0.209 0.261	0.237 0.312	0.237 0.296	0.276 0.336	0.598 0.544	0.658 0.589	0.416 0.435
	336	0.245 0.286	0.254 0.286	0.280 0.306	0.298 0.353	0.283 0.335	0.339 0.380	0.578 0.521	0.797 0.652	0.455 0.454
	720	0.324 0.342	0.326 0.337	0.365 0.359	0.352 0.288	0.345 0.381	0.403 0.428	1.059 0.741	0.869 0.675	0.535 0.520
	Avg	0.229 0.271	0.237 0.270	0.236 0.287	0.271 0.334	0.265 0.317	0.309 0.360	0.634 0.548	0.696 0.602	0.444 0.454
1 st count	24	1.974 0.886	2.063 0.881	2.317 0.934	2.527 1.000	2.398 1.040	3.228 1.260	5.764 1.677	4.480 1.444	5.914 1.734
	36	2.028 0.976	1.868 0.892	1.972 0.900	2.615 1.007	2.646 1.088	2.679 1.080	4.755 1.467	4.799 1.467	6.631 1.845
	48	2.353 1.115	1.790 0.884	2.238 0.900	2.359 0.972	2.614 1.086	2.622 1.078	4.763 1.469	4.800 1.468	6.736 1.857
	60	2.425 1.203	1.979 0.957	2.027 0.928	2.487 1.016	2.804 1.146	2.857 1.15	5.264 1.564	5.278 1.560	6.870 1.879
	Avg	2.195 1.045	1.925 0.903	2.139 0.901	2.497 1.004	2.616 1.090	2.847 1.144	5.137 1.544	4.839 1.485	6.538 1.829
1 st count		5	5	4	0	0	0	0	0	0

Table S6: Long-term Forecasting Results (MSE, MAE). TEST uses GPT2-Medium as the backbone. The past sequence length is set as 36 for ILI and 96 for the others. All the results are averaged from 4 different prediction lengths, that is $\{24, 36, 48, 60\}$ for ILI and $\{96, 192, 336, 720\}$ for the others.

2.2.4 FEW-SHOT FORECASTING

For the few-shot forecasting task, only 10% percentage timesteps of training data are used, and the other two parts remain unchanged. We follow the classical experiment settings and the results of SOTA models in Zhou et al. (2023) (NeurIPS 2023). The results are shown in Table S7. Overall, TEST has comparable performance with the SOTA baselines PatchTST and Dlinear, and SOTA LLM for TS method GPT4TS.

2.2.5 ZERO-SHOT FORECASTING

Zero-shot Forecasting task can evaluate the cross datasets adaption ability. Which means that the method is evaluated to perform on a dataset (without any training data from this dataset) when it is trained from another dataset. The results are summarized in Table S8. TEST outperforms all recent SOTA methods. TEST is comparable to N-BEATS without any meta-learning design and GPT4TS.

2.3 CLASSIFICATION TASKS

All the deep learning networks are implemented in PyTorch and trained on NVIDIA V100 32GB GPUs. We use Area Under Curve of Receiver Operating Characteristic (AUC-ROC) as metrics. Meanwhile, we compute the average rank, the number of Top-1, Top-3, and Top-5 accuracy to show the robustness of different methods. All experiments are repeated 3 times and the mean of the metrics is used in the final results.

Methods		TEST	GPT4TS	DLinear	PatchTST	TimesNet	FEDformer	Autoformer	Stationary	ETSformer	LightTS	Informer	Reformer
Weather	96	0.163 0.213	0.163 0.215	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	0.230 0.263	0.210 0.254	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	0.278 0.282	0.256 0.292	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	0.301 0.328	0.321 0.339	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	0.243 0.272	0.238 0.275	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
ETTh1	96	0.455 0.457	0.458 0.456	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	0.572 0.519	0.570 0.516	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	0.611 0.531	0.608 0.535	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	0.723 0.594	0.725 0.591	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	0.479 0.525	0.590 0.525	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
ETTh2	96	0.332 0.374	0.331 0.374	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	0.401 0.433	0.402 0.411	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	0.408 0.440	0.406 0.433	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	0.459 0.480	0.449 0.464	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	0.401 0.432	0.397 0.421	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
ETTm1	96	0.392 0.401	0.390 0.404	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	0.423 0.426	0.429 0.423	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	0.471 0.444	0.469 0.439	0.419 0.434	0.476 0.454	0.725 0.568	0.698 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	0.552 0.501	0.569 0.498	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	0.574 0.443	0.464 0.441	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
ETTm2	96	0.233 0.262	0.188 0.269	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	0.303 0.302	0.251 0.309	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	0.359 0.341	0.307 0.346	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	0.452 0.419	0.426 0.417	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	0.317 0.309	0.293 0.335	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
Electricity	96	0.143 0.235	0.139 0.237	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	0.158 0.255	0.156 0.252	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	0.176 0.275	0.175 0.270	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	0.230 0.311	0.233 0.317	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 0.269	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
Traffic	96	0.415 0.317	0.414 0.297	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	0.425 0.300	0.426 0.301	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	0.436 0.310	0.434 0.303	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	0.489 0.338	0.487 0.337	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	0.441 0.316	0.440 0.310	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
1^{st} count		5	5	4	0	0	0	0	0	0	0	0	0

Table S7: Few-shot Forecasting Results (MSE, MAE). TEST uses GPT2-Medium as the backbone. All the results are averaged from 4 different prediction lengths, that is $\{96, 192, 336, 720\}$.

Methods	M4	M3	TOURISM	ELECTR		
Metric	sMAPE	sMAPE	MAPE	NDx100	Average	1^{st} count
N-BEATS	11.70	12.44	18.82	17.8	15.19	2
DLinear	15.33	14.03	28.51	17.6	18.86	0
TimesNet	13.55	14.17	28.84	19.3	18.96	0
PatchTST	13.22	13.06	27.10	17.3	17.67	0
ETSformer	27.74	16.03	180.40	44.2	67.09	0
LightTS	13.62	17.90	66.99	19.6	29.52	0
Stationary	13.32	15.29	43.75	22.0	23.59	0
FEDformer	15.04	13.53	31.55	18.4	19.63	0
Autoformer	20.02	15.87	40.39	33.9	27.54	0
Informer	19.04	15.82	35.82	21.2	22.97	0
Reformer	14.09	13.37	25.48	21.6	18.63	0
GPT2(6)	13.12	13.06	22.14	17.2	16.38	1
TEST	13.10	12.56	18.17	17.9	15.93	1

Table S8: Zero-shot learning results. Dataset-specific metrics aggregated over each dataset. A lower value indicates better performance. The source dataset of M3, Tourism, Electricity are M4. For M4, the source data for N-BEATS is FRED, and M3 for other models.

2.3.1 DATASET DETAILS

We present accuracy scores for all 30 kinds of multivariate TS datasets in UEA archive Bagnall et al. (2018). UEA consists of 30 different datasets. Details of these datasets are shown in Table S9

Dataset	Train Cases	Test Cases	Dimensions	Length	Classes
ArticularyWordRecognition	275	30	9	144	25
AtrialFibrillation	15	15	2	640	3
BasicMotions	40	40	4	100	4
CharacterTrajectories	1422	1436	3	182	20
Cricket	108	72	6	17984	5
DuckDuckGeese	60	40	1345	270	5
EigenWorms	128	131	6	17984	5
Epilepsy	137	138	3	206	4
EthanolConcentration	261	263	3	1751	4
ERing	30	20	4	65	6
FaceDetection	5890	3524	144	62	2
FingerMovements	316	100	28	50	2
HandMovementDirection	320	147	10	400	4
Handwriting	150	850	3	152	26
Heartbeat	204	105	61	495	2
JapaneseVowels	270	370	12	29	9
Libras	180	280	2	45	15
LSST	2459	2466	6	36	14
InsectWingbeat	30000	20000	200	78	10
MotorImagery	278	100	64	3000	2
NATOPS	180	180	24	51	6
PenDigits	7494	3498	2	8	10
PEMS-SF	267	173	963	144	7
Phoneme	3315	3353	11	217	39
RacketSports	151	152	6	30	4
SelfRegulationSCP1	268	293	6	896	2
SelfRegulationSCP2	200	180	7	1152	2
SpokenArabicDigits	6599	2199	13	93	10
StandWalkJump	12	15	4	2500	3
UWaveGestureLibrary	120	320	3	315	8

Table S9: UEA Classification Dataset Details

2.3.2 BASELINE DETAILS

For classification, we refer to the SOTA methods: Three benchmarks Bostrom et al. (2018) (EDI, DTWI, and DTWD) are based on Euclidean Distance, dimension-independent dynamic time warping, and dimension-dependent dynamic time warping; MLSTM-FCNs Karim et al. (2019) applies an LSTM layer and stacked CNN layers to generate features; WEASEL-MUSE Schäfer & Leser (2017) is a bag-of-pattern based approach which extracts and represents features to words. Scalable Representation Learning (SRL) Franceschi et al. (2019a) employs negative sampling techniques with an encoder-based architecture to learn the representation; TapNet Zhang et al. (2020) is a recent model with an attentional prototype learning in its deep learning-based network; ShapeNet Li et al. (2021a) projects the subsequences into a unified space and applies clustering to find the shapelets; Rocket and MiniRocket Dempster et al. (2021) use random convolutional kernels to extract features from univariate time series; RL-PAM Gao et al. (2022) introduces reinforcement learning to the pattern mining; TStamp Transformer Zerveas et al. (2021) takes the values at each timestamp as the input for a transformer encoder; SVP-T Zuo et al. (2023) uses different variables and positions (time interval) as the inputs (shape-level).

2.3.3 MULTIVARIATE TIME SERIES CLASSIFICATION

We follow the classical experiment settings in multivariate time series classification tasks Bostrom et al. (2018). The results are shown in Table S10. Overall, TEST achieves comparable performance to SOTA models and outperforms most baselines.

2.4 REPRESENTATION TASKS

We assess the quality of our learned representations on supervised tasks in a standard manner by using them for time series classification Franceschi et al. (2019b). All the deep learning networks are implemented in PyTorch and trained on NVIDIA V100 32GB GPUs. We use Area Under Curve of Receiver Operating Characteristic (AUC-ROC) as metrics.

	EDI	DTWI	DTWD	MLSTM-FCNs	WEASEL+MUSE	SRL	TapNet	ShapeNet	Rocket	MiniRocket	RLPAM	TStamp	SVP-T	TEST
AWR	0.970	0.980	0.987	0.973	0.990	0.987	0.987	0.987	0.996	0.992	0.923	0.983	0.993	0.994
AF	0.267	0.267	0.220	0.267	0.333	0.133	0.333	0.400	0.249	0.133	0.733	0.200	0.400	0.420
BM	0.676	1.000	0.975	0.950	1.000	1.000	1.000	1.000	0.990	1.000	1.000	0.975	1.000	1.000
CT	0.964	0.969	0.989	0.985	0.990	0.994	0.997	0.980	N/A	0.993	0.978	N/A	0.990	0.989
CK	0.944	0.986	1.000	0.917	1.000	0.986	0.958	0.986	1.000	0.986	1.000	0.958	1.000	1.000
DDG	0.275	0.550	0.600	0.675	0.575	0.675	0.575	0.725	0.461	0.650	0.700	0.480	0.700	0.675
EW	0.549	N/A	0.618	0.504	0.890	0.878	0.489	0.878	0.863	0.962	0.908	N/A	0.923	0.878
EP	0.666	0.978	0.964	0.761	1.000	0.957	0.971	0.987	0.991	1.000	0.978	0.920	0.986	0.985
ER	0.133	0.914	0.929	0.133	0.133	0.133	0.133	0.133	0.981	0.981	0.819	0.933	0.937	0.937
EC	0.293	0.304	0.323	0.373	0.430	0.236	0.323	0.312	0.447	0.468	0.369	0.337	0.331	0.373
FD	0.519	0.000	0.529	0.545	0.545	0.528	0.556	0.602	0.694	0.620	0.621	0.681	0.512	0.512
FM	0.550	0.520	0.530	0.580	0.490	0.540	0.530	0.580	0.553	0.550	0.640	0.776	0.600	0.770
HMD	0.278	0.306	0.231	0.365	0.365	0.270	0.378	0.338	0.446	0.392	0.635	0.608	0.392	0.444
HW	0.200	0.316	0.286	0.286	0.605	0.533	0.357	0.452	0.567	0.507	0.522	0.305	0.433	0.431
HB	0.619	0.658	0.717	0.663	0.727	0.737	0.751	0.756	0.718	0.771	0.779	0.712	0.790	0.791
IW	0.128	N/A	N/A	0.167	N/A	0.160	0.208	0.250	N/A	0.595	0.352	0.684	0.184	0.572
JV	0.924	0.959	0.949	0.976	0.973	0.989	0.965	0.984	0.965	0.989	0.935	0.994	0.978	0.991
LB	0.833	0.894	0.870	0.856	0.878	0.867	0.850	0.856	0.906	0.922	0.794	0.844	0.883	0.884
LSST	0.456	0.575	0.551	0.373	0.590	0.558	0.568	0.590	0.632	0.643	0.643	0.381	0.666	0.595
MI	0.510	N/A	0.500	0.510	0.500	0.540	0.590	0.610	0.531	0.550	0.610	N/A	0.650	0.650
NT	0.850	0.850	0.883	0.889	0.870	0.944	0.939	0.883	0.885	0.928	0.950	0.900	0.906	0.902
PD	0.705	0.939	0.977	0.978	0.948	0.983	0.980	0.977	0.996	N/A	0.982	0.974	0.983	0.979
PM	0.973	0.734	0.711	0.699	0.000	0.688	0.751	0.751	0.856	0.522	0.632	0.919	0.867	0.860
PH	0.104	0.151	0.151	0.110	0.190	0.246	0.175	0.298	0.284	0.292	0.175	0.088	0.176	0.196
RS	0.868	0.842	0.803	0.803	0.934	0.862	0.868	0.882	0.928	0.868	0.868	0.829	0.842	0.851
SCP1	0.771	0.765	0.775	0.874	0.710	0.846	0.652	0.782	0.866	0.925	0.802	0.925	0.884	0.870
SCP2	0.483	0.533	0.539	0.472	0.460	0.556	0.550	0.578	0.514	0.522	0.632	0.589	0.600	0.579
SAD	0.967	0.959	0.963	0.990	0.982	0.956	0.983	0.975	0.630	0.620	0.621	0.993	0.986	0.982
SWJ	0.200	0.333	0.200	0.067	0.333	0.400	0.400	0.533	0.456	0.333	0.667	0.267	0.467	0.468
UGL	0.881	0.868	0.903	0.891	0.916	0.884	0.894	0.906	0.944	0.938	0.944	0.903	0.941	0.933
Avg.Rank	10.933	9.480	8.821	8.756	6.890	7.120	6.956	5.523	5.423	5.013	5.059	7.484	4.032	4.012
Num.Top-1	1	1	1	0	5	1	2	3	5	5	6	4	4	6
Num.Top-3	1	2	1	1	6	6	3	7	12	14	16	9	17	18
Num.Top-5	2	2	3	5	15	12	13	17	16	20	19	10	23	24
P-value	0.000	0.000	0.000	0.000	0.006	0.003	0.000	0.118	0.217	0.765	0.967	0.047	0.044	0.040

Table S10: Accuracies on All Datasets of the UEA Archive

2.4.1 DATASET DETAILS

We represent the results for all 128 kinds of univariate TS datasets in UCR archive Dau et al. (2019), which is a standard set of varied univariate datasets.

2.4.2 BASELINE DETAILS

The compared method includes SOTAs of unsupervised time series representation: T-Loss Franceschi et al. (2019b), TS-TCC Eldele et al. (2021), TST Zerveas et al. (2021) and TNC Tonekaboni et al. (2021), TS2Vec Yue et al. (2022).

2.4.3 CLASSIFICATION BASED ON REPRESENTATION

We assess the quality of our learned representations on supervised tasks in a standard manner by using them for time series classification Franceschi et al. (2019b). In this setting, we show that our method outperforms SOTA unsupervised methods, and notably achieves performance close to the supervised SOTA method as shown in Table S11.

For each considered dataset with a train / test split, we unsupervisedly train an encoder using its train set. We then train an SVM with radial basis function kernel on top of the learned features using the train labels of the dataset, and output the corresponding classification score on the test set.

	TEST	TCN	TS2Vec	T-Loss	TNC
Adiac	0.776	0.768	0.765	0.675	0.726
ArrowHead	0.825	0.857	0.817	0.766	0.703
Beef	0.766	0.768	0.633	0.667	0.733
BeetleFly	0.853	0.900	0.900	0.800	0.850
BirdChicken	0.808	0.803	0.800	0.850	0.750
Car	0.883	0.834	0.700	0.833	0.683
CBF	1.000	1.000	1.000	0.983	0.983

ChlorineConcentration	0.810	0.832	0.812	0.749	0.760
CinCECGTorso	0.815	0.829	0.825	0.713	0.669
Coffee	1.000	1.000	1.000	1.000	1.000
Computers	0.632	0.660	0.660	0.664	0.684
CricketX	0.802	0.787	0.805	0.713	0.623
CricketY	0.754	0.749	0.769	0.728	0.597
CricketZ	0.787	0.794	0.790	0.708	0.682
DiatomSizeReduction	0.980	0.985	0.987	0.984	0.993
DistalPhalanxOutlineCorrect	0.776	0.761	0.757	0.775	0.754
DistalPhalanxOutlineAgeGroup	0.714	0.727	0.719	0.727	0.741
DistalPhalanxTW	0.662	0.698	0.683	0.676	0.669
Earthquakes	0.746	0.748	0.748	0.748	0.748
ECG200	0.893	0.920	0.880	0.940	0.830
ECG5000	0.935	0.935	0.934	0.933	0.937
ECGFiveDays	1.000	1.000	1.000	1.000	0.999
ElectricDevices	0.714	0.721	0.719	0.707	0.700
FaceAll	0.789	0.771	0.805	0.786	0.766
FaceFour	0.834	0.932	0.932	0.920	0.659
FacesUCR	0.939	0.924	0.926	0.884	0.789
FiftyWords	0.781	0.771	0.774	0.732	0.653
Fish	0.937	0.926	0.937	0.891	0.817
FordA	0.940	0.936	0.948	0.928	0.902
FordB	0.789	0.794	0.807	0.793	0.733
GunPoint	0.983	0.980	0.987	0.980	0.967
Ham	0.714	0.714	0.724	0.724	0.752
HandOutlines	0.918	0.925	0.930	0.922	0.930
Haptics	0.510	0.526	0.536	0.490	0.474
Herring	0.625	0.644	0.609	0.594	0.594
InlineSkate	0.389	0.418	0.407	0.371	0.378
InsectWingbeatSound	0.620	0.630	0.624	0.597	0.549
ItalyPowerDemand	0.969	0.925	0.960	0.954	0.928
LargeKitchenAppliances0	0.855	0.845	0.875	0.789	0.776
Lightning2	0.846	0.869	0.820	0.869	0.869
Lightning7	0.866	0.863	0.822	0.795	0.767
Mallat	0.915	0.944	0.873	0.951	0.871
Meat	0.950	0.952	0.967	0.950	0.917
MedicalImages	0.792	0.789	0.793	0.750	0.754
MiddlePhalanxOutlineCorrect	0.811	0.838	0.825	0.825	0.818
MiddlePhalanxOutlineAgeGroup	0.636	0.636	0.630	0.656	0.643
MiddlePhalanxTW	0.591	0.584	0.578	0.591	0.571
MoteStrain	0.857	0.861	0.863	0.851	0.825
NonInvasiveFetalECGThorax1	0.923	0.930	0.919	0.878	0.898
NonInvasiveFetalECGThorax2	0.940	0.938	0.935	0.919	0.912
OliveOil	0.903	0.901	0.940	0.867	0.833
OSULeaf	0.872	0.851	0.843	0.760	0.723
PhalangesOutlinesCorrect	0.794	0.809	0.823	0.784	0.787
Phoneme	0.296	0.312	0.309	0.276	0.180
Plane	1.000	1.000	0.990	0.990	1.000
ProximalPhalanxOutlineCorrect	0.876	0.887	0.900	0.859	0.866
ProximalPhalanxOutlineAgeGroup	0.844	0.837	0.829	0.844	0.854
ProximalPhalanxTW	0.785	0.824	0.805	0.771	0.810
RefrigerationDevices	0.587	0.586	0.589	0.515	0.565
ScreenType	0.405	0.414	0.397	0.416	0.509
ShapeletSim	0.989	1.000	0.994	0.672	0.589
ShapesAll	0.897	0.902	0.905	0.848	0.788
SmallKitchenAppliances	0.723	0.731	0.733	0.677	0.725
SonyAIBORobotSurface1	0.874	0.903	0.900	0.902	0.804
SonyAIBORobotSurface2	0.893	0.871	0.889	0.889	0.834
StarLightCurves	0.970	0.968	0.971	0.964	0.968
Strawberry	0.962	0.966	0.965	0.954	0.951
SwedishLeaf	0.939	0.945	0.942	0.914	0.880
Symbols	0.973	0.977	0.972	0.963	0.885
SyntheticControl	0.997	0.997	0.993	0.987	1.000
ToeSegmentation1	0.933	0.917	0.947	0.939	0.864
ToeSegmentation2	0.915	0.899	0.900	0.900	0.831
Trace	1.000	1.000	1.000	0.990	1.000
TwoLeadECG	0.982	0.986	0.987	0.999	0.993
TwoPatterns	1.000	1.000	1.000	0.999	1.000
UWaveGestureLibraryX	0.810	0.795	0.801	0.785	0.781
UWaveGestureLibraryY	0.729	0.719	0.720	0.710	0.697
UWaveGestureLibraryZ	0.761	0.774	0.768	0.757	0.721
UWaveGestureLibraryAll	0.935	0.930	0.934	0.896	0.903
Wafer	0.995	0.998	0.998	0.992	0.994
Wine	0.788	0.880	0.889	0.815	0.759
WordSynonyms	0.699	0.679	0.704	0.691	0.630
Worms	0.704	0.701	0.701	0.727	0.623
WormsTwoClass	0.805	0.806	0.753	0.792	0.727
Yoga	0.883	0.883	0.877	0.837	0.812
ACSF1	0.849	0.910	0.910	0.900	0.730
AllGestureWiimoteX	0.744	0.777	0.751	0.763	0.703
AllGestureWiimoteY	0.754	0.796	0.774	0.726	0.699

AllGestureWiimoteZ	0.744	0.749	0.770	0.723	0.646
BME	0.979	0.992	0.980	0.993	0.973
Chinatown	0.969	0.964	0.959	0.951	0.977
Crop	0.753	0.754	0.758	0.722	0.738
EOGHorizontalSignal	0.544	0.569	0.522	0.605	0.442
EOGVerticalSignal	0.467	0.503	0.472	0.434	0.392
EthanolLevel	0.480	0.468	0.484	0.382	0.424
FreezerRegularTrain	0.983	0.996	0.983	0.956	0.991
FreezerSmallTrain	0.893	0.875	0.872	0.933	0.982
Fungi	0.967	0.958	0.946	1.000	0.527
GestureMidAirD1	0.637	0.608	0.615	0.608	0.431
GestureMidAirD2	0.508	0.479	0.515	0.546	0.362
GestureMidAirD3	0.346	0.492	0.300	0.285	0.292
GesturePebbleZ1	0.878	0.930	0.884	0.919	0.378
GesturePebbleZ2	0.842	0.873	0.848	0.899	0.316
GunPointAgeSpan	0.994	0.987	0.968	0.994	0.984
GunPointMaleVersusFemale	1.000	1.000	1.000	0.997	0.994
GunPointOldVersusYoung	1.000	1.000	1.000	1.000	1.000
HouseTwenty	0.944	0.917	0.941	0.933	0.782
InsectEPGRegularTrain	1.000	1.000	1.000	1.000	1.000
InsectEPGSmallTrain	1.000	1.000	1.000	1.000	1.000
MelbournePedestrian	0.954	0.959	0.956	0.944	0.942
MixedShapesRegularTrain	0.915	0.917	0.922	0.905	0.911
MixedShapesSmallTrain	0.884	0.861	0.856	0.860	0.813
PickupGestureWiimoteZ	0.800	0.823	0.760	0.740	0.620
PigAirwayPressure	0.524	0.630	0.683	0.510	0.413
PigArtPressure	0.962	0.966	0.966	0.928	0.808
PigCVP	0.803	0.815	0.870	0.788	0.649
PLAID	0.551	0.561	0.549	0.555	0.495
PowerCons	0.967	0.961	0.972	0.900	0.933
Rock	0.660	0.700	0.700	0.580	0.580
SemgHandGenderCh2	0.952	0.963	0.962	0.890	0.882
SemgHandSubjectCh2	0.897	0.860	0.891	0.789	0.593
SemgHandMovementCh2	0.944	0.952	0.942	0.920	0.820
SmoothSubspace	0.967	0.980	0.993	0.960	0.913
UMD	1.000	1.000	0.993	0.993	0.993
Avg	0.826	0.832	0.827	0.806	0.761

Table S11: Accuracies on All Datasets of the UCR Archive

2.5 ABLATION

TEST contains two contrastive learning strategies: instance-wise contrast and feature-wise contrast, and can use different text embedding vectors as prototypes, we show the impact of these strategies.

2.5.1 CONTRASTIVE LEARNING STRATEGIES

As shown in Table S12 and S13, both two contrastive learning strategies can increase the accuracy.

	ETTm1	ETTm2	ETTh1	ETTh2	Electricity	Traffic	Weather	ILI
Instance-wise	0.621 0.550	0.755 0.630	0.493 0.453	0.580 0.612	0.293 0.396	0.788 0.620	0.463 0.349	3.301 4.535
Feature-wise	0.741 0.559	0.793 0.634	0.699 0.493	0.585 0.628	0.286 0.390	0.821 0.629	0.453 0.388	3.139 5.931
TEST	0.353 0.382	0.293 0.334	0.414 0.431	0.331 0.380	0.162 0.253	0.430 0.295	0.229 0.271	2.195 1.045

Table S12: Long-term Forecasting Results (MSE, MAE). TEST uses different contrastive learning strategy. All the results are averaged from 4 different prediction lengths, that is $\{24, 36, 48, 60\}$ for ILI and $\{96, 192, 336, 720\}$ for the others. The results are average.

	TEST	Instance-wise	Feature-wise	TimesNet	N-BEATS	ETSformer	DLinear	FEDformer	Stationary	Autoformer	Informer	Reformer
SMAPE	11.927	13.525	16.987	11.829	11.851	14.718	13.639	12.840	12.780	12.909	14.086	18.200
MASE	1.613	2.111	3.265	1.585	1.599	2.408	2.095	1.701	1.756	1.771	3.010	4.223
OWA	0.861	1.051	1.480	0.851	0.855	1.172	1.051	0.918	0.930	0.939	1.230	1.775

Table S13: Short-term Forecasting Task on M4. The prediction lengths are in $[6, 48]$ and results are averaged from several datasets.

2.5.2 TEXT PROTOTYPES

The number and the type of text prototypes will lead to different results.

As shown in Table S14. We randomly select 1, 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22 prototypes. The accuracy and number are basically positively correlated. The results of 10 prototypes are almost optimal.

As shown in Table S15. We randomly select 10 prototypes 10 times. The accuracy is basically consistent. Therefore, the type of prototypes has almost no impact on the results.

	1	2	4	6	8	10	12	14	16	18	20	22
SMAPE	30.901	20.201	17.415	16.997	13.820	11.927	11.710	11.638	11.094	11.098	10.953	10.885
MASE	6.590	4.515	3.910	3.595	2.580	1.613	1.408	1.195	1.301	1.306	1.471	1.310
OWA	3.779	2.050	1.451	1.484	0.990	0.861	0.872	0.801	0.910	0.902	0.838	0.830

Table S14: Short-term Forecasting Task on M4. The results are reported with different number of text prototypes.

	1	2	3	4	5	6	7	8	9	10	Avg.	Std.
SMAPE	11.907	11.920	11.927	11.926	11.925	11.925	11.950	11.890	11.728	11.910	11.901	0.059
MASE	1.612	1.610	1.653	1.603	1.619	1.620	1.625	1.623	1.613	1.591	1.617	0.016
OWA	0.870	0.872	0.872	0.872	0.872	0.872	0.849	0.862	0.876	0.870	0.868	0.009

Table S15: Short-term Forecasting Task on M4. The results are reported with different types of text prototypes.

Considering why the type of text prototype does not significantly affect results, we figure that in high dimensional space, almost all vectors are pairwise orthogonal Hopcroft & Kannan (2013). Which means that, in high-dimensional space, it is easy to generate a large number of almost orthogonal vectors to represent different attributes. Thus, randomly selecting the same number of vectors, the represented space size and expressed number of features are almost the same. Therefore, the key is the number rather than the type.

In terms of probability, “two vectors orthogonal” is equivalent to “two vectors perpendicular” is equivalent to “two vectors uncorrelated” is equivalent to “ $\cos \theta = 0$ ”. For a n -dimensional space, randomly two vectors have: $\forall \epsilon, \lim_{n \rightarrow \infty} P(|\cos \theta| > \epsilon) = 0$. As shown in Figure S2, as the dimension increases, the probability of two random vectors being similar decreases. For LLM, $n > 1024, P(\theta = 0) < 0.00001$.

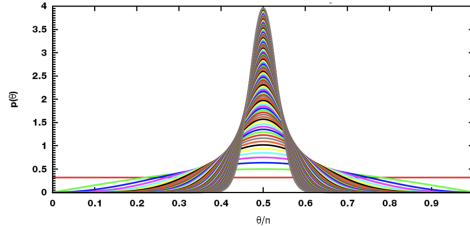


Figure S2: Probability Density of the Angle between Two Random Vectors in n-dimensional Space

REFERENCES

- Anthony J. Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul Southam, and Eamonn J. Keogh. The UEA multivariate time series classification archive, 2018. *CoRR*, abs/1811.00075, 2018.
- Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, pp. 1476–4687, 2023. doi: 10.1038/s41586-023-06545-z.
- Aaron Bostrom, Anthony Bagnall, Eamonn Keogh, Hoang Anh Dau, James Large, Jason Lines, Michael Flynn, and Paul Southam. The uea multivariate time series classification archive, 2018, 2018.

- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. 2020.
- CDC. Illness. 2021. doi: <https://gis.cdc.gov/grasp/fluview/fluportaldashboard.html>.
- Ching Chang, Wen-Chih Peng, and Tien-Fu Chen. Llm4ts: Two-stage fine-tuning for time-series forecasting with pre-trained llms, 2023.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of International Conference on Machine Learning*, volume 119, pp. 1597–1607, 2020.
- Hyunseung Chung, Jiho Kim, Joon myoung Kwon, Ki-Hyun Jeon, Min Sung Lee, and Edward Choi. Text-to-ecg: 12-lead electrocardiogram synthesis conditioned on clinical text reports, 2023.
- Hoang Anh Dau, Anthony Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn Keogh. The ucr time series archive. *IEEE/CAA Journal of Automatica Sinica*, 6:1293–1305, 2019. doi: 10.1109/JAS.2019.1911747.
- Angus Dempster, Daniel F. Schmidt, and Geoffrey I. Webb. Minirocket: A very fast (almost) deterministic transform for time series classification. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 248–257, 2021. doi: 10.1145/3447548.3467231.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of Annual Meeting of the Association for Computational Linguistics*, volume 1, pp. 320–335, 2022.
- Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee Keong Kwoh, Xiaoli Li, and Cuntai Guan. Time-series representation learning via temporal and contextual contrasting. In *International Joint Conference on Artificial Intelligence*, pp. 2352–2359, 2021.
- Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised scalable representation learning for multivariate time series. In *Advances in Neural Information Processing Systems*, pp. 4652–4663, 2019a.
- Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised scalable representation learning for multivariate time series. In *Advances in Neural Information Processing Systems*, pp. 4652–4663, 2019b.
- Ge Gao, Qitong Gao, Xi Yang, Miroslav Pajic, and Min Chi. A reinforcement learning-informed pattern mining framework for multivariate time series classification. In *Proceedings of International Joint Conference on Artificial Intelligence*, pp. 2994–3000, 2022. doi: 10.24963/IJCAI.2022/415.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent - A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, 2020.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew Gordon Wilson. Large language models are zero-shot time series forecasters. *CoRR*, abs/2310.07820, 2023. doi: 10.48550/ARXIV.2310.07820.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *Computer Vision and Pattern Recognition*, pp. 9726–9735, 2020.
- John Hopcroft and Ravindran Kannan. *Computer science theory for the information age*. Cambridge University press, 2013.

- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y. Zhang, Xiaoming Shi, Pin-Yu Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. Time-llm: Time series forecasting by reprogramming large language models. *CoRR*, abs/2310.01728, 2023. doi: 10.48550/ARXIV.2310.01728.
- Fazle Karim, Somshubra Majumdar, Houshang Darabi, and Samuel Harford. Multivariate lstm-fcns for time series classification. *Neural Networks*, 116:237–245, 2019. doi: 10.1016/J.NEUNET.2019.04.014.
- Salar Hosseini Khorasgani, Yuxuan Chen, and Florian Shkurti. SLIC: self-supervised learning with iterative clustering for human action videos. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16070–16080, 2022. doi: 10.1109/CVPR52688.2022.01562.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.
- Guozhong Li, Byron Choi, Jianliang Xu, Sourav S. Bhowmick, Kwok-Pan Chun, and Grace Lai-Hung Wong. Shapenet: A shapelet-neural network approach for multivariate time series classification. In *AAAI Conference on Artificial Intelligence*, pp. 8375–8383, 2021a. doi: 10.1609/AAAI.V35I9.17018.
- Jun Li, Che Liu, Sibbo Cheng, Rossella Arcucci, and Shenda Hong. Frozen language model helps ECG zero-shot learning. *CoRR*, abs/2303.12311, 2023. doi: 10.48550/arXiv.2303.12311.
- Junnan Li, Pan Zhou, Caiming Xiong, and Steven C. H. Hoi. Prototypical contrastive learning of unsupervised representations. In *International Conference on Learning Representations*, 2021b.
- Yunfan Li, Peng Hu, Jerry Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering. In *AAAI Conference on Artificial Intelligence*, pp. 8547–8555, 2021c.
- Xin Liu, Daniel McDuff, Geza Kovacs, Isaac R. Galatzer-Levy, Jacob E. Sunshine, Jiening Zhan, Ming-Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak N. Patel. Large language models are few-shot health learners. *CoRR*, abs/2305.15525, 2023. doi: 10.48550/arXiv.2305.15525.
- Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. In *Advances in Neural Information Processing Systems*, 2022.
- Qianli Ma, Zhen Liu, Zhenjing Zheng, Ziyang Huang, Siying Zhu, Zhongzhong Yu, and James T. Kwok. A survey on time-series pre-trained models. *CoRR*, abs/2305.10716, 2023. doi: 10.48550/arXiv.2305.10716.
- Qianwen Meng, Hangwei Qian, Yong Liu, Yonghui Xu, Zhiqi Shen, and Lizhen Cui. MHCCCL: masked hierarchical cluster-wise contrastive learning for multivariate time series. *CoRR*, abs/2212.01141, 2022.
- Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*, 2023.
- Boris N. Oreshkin, Dmitri Carpo, Nicolas Chapados, and Yoshua Bengio. N-BEATS: neural basis expansion analysis for interpretable time series forecasting. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- PeMS. Traffic. 2021. doi: <http://pems.dot.ca.gov/>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019.
- Patrick Schäfer and Ulf Leser. Multivariate time series classification with WEASEL+MUSE. *CoRR*, abs/1711.11343, 2017.
- Vivek Sharma, Makarand Tapaswi, M. Saquib Sarfraz, and Rainer Stiefelhagen. Clustering based contrastive learning for improving face representations. In *IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 109–116, 2020. doi: 10.1109/FG47880.2020.00011.

- Taylor SJ and Letham B. Forecasting at scale. In *PeerJ Preprints*, pp. 5:e3190v2, 2017. doi: 10.7287/peerj.preprints.3190v2.
- Sana Tonekaboni, Danny Eytan, and Anna Goldenberg. Unsupervised representation learning for time series with temporal neighborhood coding. In *International Conference on Learning Representations*, 2021.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, and et al. Llama 2: Open foundation and fine-tuned chat models, 2023.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.
- Wetterstation. Weather. 2017. doi: <https://www.bgc-jena.mpg.de/wetter/>.
- Kristoffer Wickstrøm, Michael Kampffmeyer, Karl Øyvind Mikalsen, and Robert Jenssen. Mixing up contrastive learning: Self-supervised representation learning for time series. *Pattern Recognit. Lett.*, 155:54–61, 2022.
- Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven C. H. Hoi. Cost: Contrastive learning of disentangled seasonal-trend representations for time series forecasting. In *International Conference on Learning Representations*, 2022a.
- Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven C. H. Hoi. Etsformer: Exponential smoothing transformers for time-series forecasting. *CoRR*, abs/2202.01381, 2022b.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *Advances in Neural Information Processing Systems*, pp. 22419–22430, 2021.
- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *International Conference on Learning Representations*, 2023.
- Hao Xue and Flora D. Salim. Promptcast: A new prompt-based learning paradigm for time series forecasting. *CoRR*, abs/2210.08964, 2023. doi: 10.48550/arXiv.2210.08964.
- Ling Yang and Shenda Hong. Unsupervised time-series representation learning with iterative bilinear temporal-spectral fusion. In *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 25038–25054, 2022.
- Xinyu Yang, Zhenguo Zhang, and Rongyi Cui. Timeclr: A self-supervised contrastive learning framework for univariate time series representation. *Knowl. Based Syst.*, 245:108606, 2022.
- Zhihan Yue, Yujing Wang, Juanyong Duan, Tianmeng Yang, Congrui Huang, Yunhai Tong, and Bixiong Xu. Ts2vec: Towards universal representation of time series. In *AAAI Conference on Artificial Intelligence*, pp. 8980–8987, 2022.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *AAAI Conference on Artificial Intelligence*, pp. 11121–11128, 2023. doi: 10.1609/aaai.v37i9.26317.
- George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2114–2124, 2021. doi: 10.1145/3447548.3467401.
- Dejiao Zhang, Feng Nan, Xiaokai Wei, Shang-Wen Li, Henghui Zhu, Kathleen R. McKeown, Ramesh Nallapati, Andrew O. Arnold, and Bing Xiang. Supporting clustering with contrastive learning. In *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5419–5430, 2021. doi: 10.18653/V1/2021.NAACL-MAIN.427.

- Xuchao Zhang, Yifeng Gao, Jessica Lin, and Chang-Tien Lu. Tapnet: Multivariate time series classification with attentional prototypical network. In *AAAI Conference on Artificial Intelligence*, pp. 6845–6852, 2020. doi: 10.1609/AAAI.V34I04.6165.
- Xiaochen Zheng, Xingyu Chen, Manuel Schürch, Amina Mollaysa, Ahmed Allam, and Michael Krauthammer. Simts: Rethinking contrastive representation learning for time series forecasting. *CoRR*, abs/2303.18205, 2023.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI Conference on Artificial Intelligence*, pp. 11106–11115, 2021. doi: 10.1609/aaai.v35i12.17325.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 27268–27286, 2022.
- Tian Zhou, PeiSong Niu, Xue Wang, Liang Sun, and Rong Jin. One fits all: power general time series analysis by pretrained lm, 2023.
- Rundong Zuo, Guozhong Li, Byron Choi, Sourav S. Bhowmick, Daphne Ngar-yin Mah, and Grace Lai-Hung Wong. SVP-T: A shape-level variable-position transformer for multivariate time series classification. In *AAAI Conference on Artificial Intelligence*, pp. 11497–11505, 2023. doi: 10.1609/AAAI.V37I9.26359.