

CaSS: A Channel-aware Self-supervised Representation Learning Framework for Multivariate Time Series Classification

Yijiang Chen¹, Xiangdong Zhou¹, Zhen Xing¹, Zhidan Liu¹, and Minyang Xu^{1,2}

¹ School of Computer Science, Fudan University, Shanghai, China, 200433
{chenyj20,xdzhou,zxing20,zdliu20,16110240027}@fudan.edu.cn

² Arcplus Group PLC, Shanghai, China, 200041

Abstract. Self-supervised representation learning of Multivariate Time Series (MTS) is a challenging task and attracts increasing research interests in recent years. Many previous works focus on the pretext task of self-supervised learning and usually neglect the complex problem of MTS encoding, leading to unpromising results. In this paper, we tackle this challenge from two aspects: encoder and pretext task, and propose a unified channel-aware self-supervised learning framework CaSS. Specifically, we first design a new Transformer-based encoder Channel-aware Transformer (CaT) to capture the complex relationships between different time channels of MTS. Second, we combine two novel pretext tasks Next Trend Prediction (NTP) and Contextual Similarity (CS) for the self-supervised representation learning with our proposed encoder. Extensive experiments are conducted on several commonly used benchmark datasets. The experimental results show that our framework achieves new state-of-the-art comparing with previous self-supervised MTS representation learning methods (up to +7.70% improvement on LSST dataset) and can be well applied to the downstream MTS classification.

1 Introduction

With the fast progress of IoT and coming 5Gs, multivariate time series widely exists in medical, financial, industrial and other fields as an increasingly important data form [6,8]. Compared with univariate time series, multivariate time series usually contains more information and brings more potentials for data mining, knowledge discovery and decision making, etc.. However, MTS not only contains time-wise patterns, but also has complex relationships between different channels, which makes MTS analysis much more difficult.

In recent years, the self-supervised learning attracts more and more attention from research and industry communities. The self-supervised pre-training demonstrates its success in the fields of Natural Language Processing (NLP) [17] and Computer Visions (CV) [29]. Especially in NLP, adopting a pre-trained language model is de facto the first step of almost all the NLP tasks. Likewise,



Fig. 1. A sample of encoding of multivariate time series.

the self-supervised representation learning of time series not only brings performance improvements, but also helps to close the gap between the increasing amount of data (more abstraction and complexity) and the expensive cost of manually labeling for supervised learning tasks. Many research efforts have been devoted to the self-supervised representation learning of time series [15,12,34] and promising results have been achieved. However, the previous works for MTS are limited and the challenge still exists.

The self-supervised representation learning usually consists of two aspects: encoder and pretext task. As shown in Figure 1, in the past works, most of them only focus on time-wise features where all channel values of one or several time steps are fused through convolution or fully connected layer directly in the embedding process. There is a lack of deliberate investigation of the relationships between channel-wise features, which affects the encoder’s ability to capture the whole characteristics of the MTS. To deal with the problem, the methods combining with Recurrent Neural Network (RNN) are presented to capture the individual feature of each channel [3,35]. The recent work of [21] employs Transformer [30] to integrate the features of time-wise and channel-wise. Among these solutions, RNN seems not very suitable for self-supervised learning due to the consumption and is usually employed in prediction task [14]. Transformer is becoming more and more popular and is suitable for time series [34], however the previous Transformer-based MTS embedding provokes the problem of high complexity of computing and space, which prevents it from real applications. It inspires us to design a more effective Transformer-based for MTS to take advantages of the strong encoding ability of Transformer. In the aspect of pretext task, most of the previous works adopt the traditional time series embedding based on time-wise features. How to integrate channel-wise features with pretext task is a challenge.

In this paper, we propose a novel self-supervised learning framework CaSS from the aspects of encoder and pretext task. First, we propose a new Transformer-based encoder Channel-aware Transformer (CaT) for MTS encoding which investigates time-wise and channel-wise features simultaneously. It is noticed that in practice the number of channels of MTS is fixed while the time length can be unlimited and the number of channels is usually much less than the time length.

Therefore as shown in Figure 1, different from previous work [21], we integrate the time-wise features into the channel-wise features and concatenate all these novel channel-wise features as the representation of the sample. Second, we design a new self-supervised pretext task Next Trend Prediction (NTP) from the perspective of channel-wise for the first time in self-supervised MTS representation learning. It is considered that in many cases only the rise and fall of future time rather than the specific value of the time series is necessary. So we cut the multivariate time series from the middle, and using the previous sequences of all rest channels to predict the trend for each channel. Different from fitting the specific value (regression), the prediction of trend (rise and fall) is more suitable for arbitrary data. We also demonstrate through experiments that compared with fitting specific values, prediction of trend is more efficient. In addition, we employ another task called Contextual Similarity (CS) which combines a novel data augmentation strategy to maximize the similarity between similar samples and learn together with NTP task. The CS task focuses on the difference between samples while NTP task focuses on the sample itself and helps to learn the complex internal characteristics.

In summary, the main contributions of our work are as follows:

- We propose a new Transformer-based encoder Channel-aware Transformer. It can efficiently integrate the time-wise features to the channel-wise representation.
- We design two novel pretext tasks, Next Trend Prediction and Contextual Similarity for our CaSS framework. To the best of our knowledge, Next Trend Prediction task is conducted from the perspective of channel-wise for the first time in self-supervised MTS representation learning.
- We conduct extensive experiments on several commonly used benchmark datasets from different fields. Compared with the state-of-the-art self-supervised MTS representation learning methods, our method achieves new state-of-the-art in MTS classification (up to +7.70% improvement on LSST dataset). We also demonstrate its ability in few-shot learning.

2 Related Work

2.1 Encoders for Time Series Classification

A variety of methods have been proposed for time series classification. Early works employ traditional machine learning methods to solve the problem, like combining Dynamic Time Warping (DTW) [32] with Support Vector Machine (SVM) [9]. Time Series Forest [11] introduces an approach based on Random Forest [7]. Bag of Patterns (BOP) [18] and Bag of SFA Symbols (BOSS) [28] construct a dictionary-based classifier. Although these early works can deal with the problem to some extent, they need heavy crafting on data preprocessing and feature engineering. The emergence of deep learning greatly reduces feature engineering and boosts the performance of many machine learning tasks. So far Convolutional Neural Network (CNN) is popular in time series classification due

to its balance between its effect and cost, such as Multi-scale Convolutional Neural Network (MCNN) [10] for univariate time series classification, Multi-channels Deep Convolutional Neural Networks (MC-DCNN) [38] for multivariate time series classification and so on [36,20,37]. Hierarchical Attention-based Temporal Convolutional Network (HA-TCN) [19] and WaveATTentionNet (WATTNet) [25] apply dilated causal convolution to improve the encoder’s effect. Among these CNN methods, Fully Convolutional Network (FCN) and Residual Network (ResNet) [31] have been proved to be the most powerful encoders in multivariate time series classification task [14]. Due to the high computation complexity, RNN based encoders are rarely applied solely to the time series classification [22,35]. It is often combined with CNN to form a two tower structure [16,3]. In recent years, more and more works have tried to apply Transformer to time series [33,26,39]. However, most of them are designed for prediction task, and few works cover the classification problem [24,27,21,34].

2.2 Pretext Tasks for Time Series

Manually labeling is a long lasting challenge for the supervised learning, and recently self-supervised training (no manually labeling) becomes more and more popular in many research fields including time series analysis. To name a few, [15] employs the idea of word2vec [23] which regards part of the time series as word, the rest as context, and part of other time series as negative samples for training. [12] employs the idea of contrastive learning where two positive samples are generated by weak augmentation and strong augmentation to predict each other while the similarity among different augmentations of the same sample is maximized. [13] is designed for univariate time series. It samples several segments of the time series and labels each segment pairs according to their relative distance in the origin series. It also adds the task of judging whether two segments are generated by the same sample. [4] is based on sampling pairs of time windows and predicting whether time windows are close in time by setting thresholds to learn EEG features. [34] is a Transformer-based method which employs the idea of mask language model [17]. The mask operation is performed for multivariate time series and the encoder is trained by predicting the masked value. However, the previous works usually focus on time-wise features and need to continuously obtain the features of several time steps [15,12,34], which makes them difficult to be applied to the novel MTS representation.

3 The Framework

In this work, we focus on self-supervised multivariate time series representation learning. Given M multivariate time series $X = \{x_0, \dots, x_M\}$, where $x_i \in \mathbb{R}^{C \times T}$ refers to the i -th time series which has C channels and T time steps. For each multivariate time series x_i , our goal is to generate a proper representation z_i which is applicable to the subsequent task.

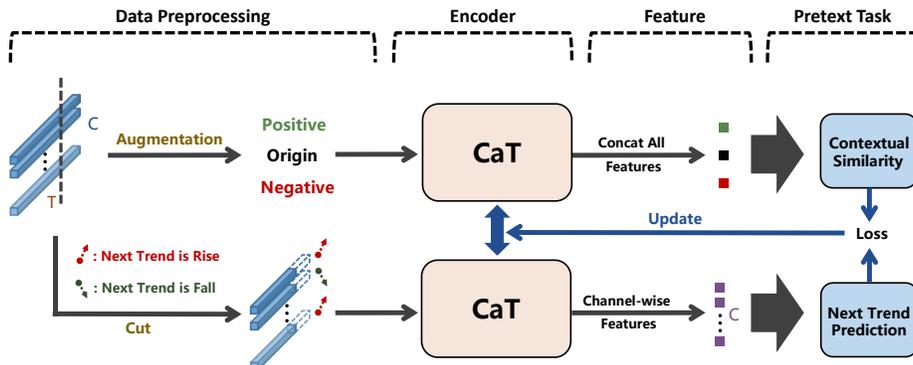


Fig. 2. Overall architecture of the channel-aware self-supervised framework CaSS.

Our proposed novel encoder Channel-aware Transformer (CaT) and pretext tasks constitute our channel-aware self-supervised learning framework CaSS. The overall architecture is shown in Figure 2. In our framework, CaT is to generate the novel channel-wise features of the MTS sample, and the generated features are served as the inputs of Next Trend Prediction task and Contextual Similarity task. Specifically, we first preprocess the time series samples and apply the two pretext tasks to learn the encoder (representations), then we employ the learnt representations to the MTS classification task by freezing the encoder.

4 Channel-aware Transformer

This section describes our proposed Transformer-based encoder Channel-aware Transformer which is served as the encoder in our self-supervised learning framework. As shown in Figure 3, It consists of Embedding Layer, Co-Transformer Layer and Aggregate Layer. The two Transformer structures in Co-Transformer Layer encode the time-wise and channel-wise features respectively, and interact with each other during encoding. Finally, we fuse the time-wise features into channel-wise features through Aggregate Layer to generate the final representation.

4.1 Embedding Layer

Given an input sample $x \in \mathbb{R}^{C \times T}$, we map it to the D -dimension time vector space and channel vector space respectively to obtain the time embedding $e_t \in \mathbb{R}^{T \times D}$ and channel embedding $e_c \in \mathbb{R}^{C \times D}$:

$$e_t = x^T W_t + e_{pos}, \quad (1)$$

$$e_c = x W_c, \quad (2)$$

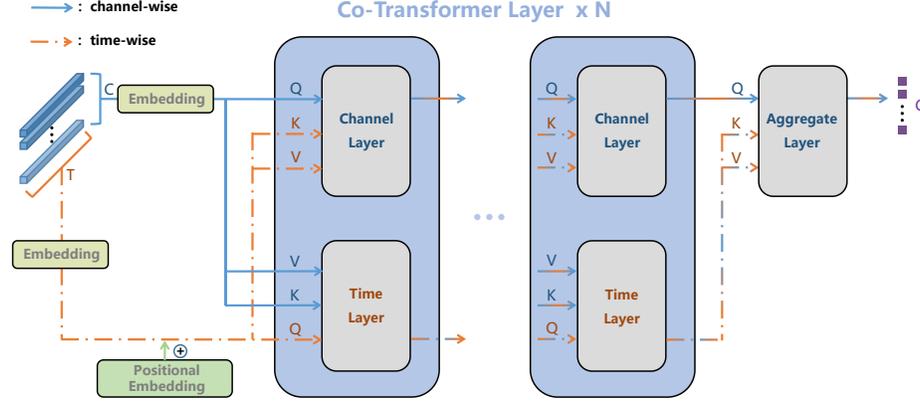


Fig. 3. Overall architecture of Channel-aware Transformer (CaT). Q, K, V represent the query vector, key vector, value vector in attention operation respectively.

where $W_t \in \mathbb{R}^{C \times D}$, $W_c \in \mathbb{R}^{T \times D}$ are learnable embedding matrices. $e_{pos} \in \mathbb{R}^{T \times D}$ is the positional embedding applying the design of [30].

4.2 Co-Transformer Layer

This part adopts a N -layer two tower structure based on Transformer [21]. Each layer is composed of Time Layer and Channel Layer, focusing on time-wise and channel-wise respectively. Supposing in the i -th layer ($i = 0, \dots, N-1$), we obtain the input $a_t^i \in \mathbb{R}^{T \times D}$ for Time Layer and $a_c^i \in \mathbb{R}^{C \times D}$ for Channel Layer. The process in Time Layer is:

$$Q_t^i = a_t^i W_{qt}^i, \quad K_t^i = a_c^i W_{kt}^i, \quad V_t^i = a_c^i W_{vt}^i, \quad (3)$$

$$b_t^i = \text{LayerNorm}(\text{MHA}(Q_t^i, K_t^i, V_t^i) + a_t^i), \quad (4)$$

$$a_t^{i+1} = \text{LayerNorm}(\text{FFN}(b_t^i) + b_t^i), \quad (5)$$

and the process in Channel Layer is:

$$Q_c^i = a_c^i W_{qc}^i, \quad K_c^i = a_t^i W_{kc}^i, \quad V_c^i = a_t^i W_{vc}^i, \quad (6)$$

$$b_c^i = \text{LayerNorm}(\text{MHA}(Q_c^i, K_c^i, V_c^i) + a_c^i), \quad (7)$$

$$a_c^{i+1} = \text{LayerNorm}(\text{FFN}(b_c^i) + b_c^i), \quad (8)$$

where $W_{qt}^i, W_{kt}^i, W_{vt}^i, W_{qc}^i, W_{kc}^i, W_{vc}^i \in \mathbb{R}^{D \times D}$ are learnable matrices. MHA is the abbreviation of multi-head attention and FFN is the abbreviation of feed forward network. Specifically, $a_t^0 = e_t, a_c^0 = e_c$.

This interactive way helps reduce the time complexity from $o((T^2 + C^2)D)$ to $o(2TCD)$ compared with the non-interactive two tower Transformer-based encoder which applies self-attention mechanism [30]. In the real world, T is usually much larger than C , so it can further boost the speed of the encoder.

4.3 Aggregate Layer

Through Co-Transformer Layer, we can obtain time-wise features a_t^N and channel-wise features a_c^N . If all features are forcibly concatenated, the final dimension of the representation will be too large to be applied in the subsequent work. In real applications, the channel length is usually much less than the time length, therefore we integrate the time-wise features into the channel-wise features through attention operation. Finally, we concatenate these novel channel-wise features as the final representation $z \in \mathbb{R}^{1 \times (C \cdot D)}$:

$$Q_c^N = a_c^N W_{qc}^N, \quad K_c^N = a_t^N W_{kc}^N, \quad V_c^N = a_t^N W_{vc}^N, \quad (9)$$

$$a_c = \text{MHA}(Q_c^N, K_c^N, V_c^N), \quad (10)$$

$$z = [a_c^1, a_c^2, \dots, a_c^C], \quad (11)$$

where $W_{qc}^N, W_{kc}^N, W_{vc}^N \in \mathbb{R}^{D \times D}$ are learnable matrices. $[\cdot, \cdot]$ is the concatenation operation.

5 Pretext Task

In order to enable our encoder to carry out self-supervised learning more efficiently, we design two novel pretext tasks Next Trend Prediction (NTP) and Contextual Similarity (CS) based on our novel channel-wise representation.

5.1 Next Trend Prediction

Given a sample $x_i \in \mathbb{R}^{C \times T}$, we randomly select a time point $t \in [1, T - 1]$ for truncation. The sequence before t time step $x_i^{NTP(t)} \in \mathbb{R}^{C \times t}$ is regarded as the input of the NTP task and the data after t is padded to T with zeros. For each channel j , we adopt the trend of the $t + 1$ time step as the label $y_{i,j}^{NTP(t)}$ for training:

$$y_{i,j}^{NTP(t)} = \begin{cases} 1, & \text{if } x_i[j, t + 1] \geq x_i[j, t] \\ 0, & \text{if } x_i[j, t + 1] < x_i[j, t] \end{cases}, \quad (12)$$

where $x_i[j, t]$ represents the value of t time step of channel j in x_i .

After inputting the NTP sample x_i into our encoder, we can obtain the representation $z_i^{NTP(t)} \in \mathbb{R}^{C \times D}$ where $z_{i,j}^{NTP(t)} \in \mathbb{R}^D$ represents the representation of the j -th channel. Finally, a projection head is applied to predict the probability

of the rise and fall. Assuming that every sample generated K_{NTP} input samples and the corresponding truncating time point set is $S \in \mathbb{R}^{K_{NTP}}$. For x_i , the loss of the NTP task can be obtained through the following formula:

$$\ell_{NTP} = \sum_{t \in S} \sum_j^C \text{CE}(\varphi_0(z_{i,j}^{NTP(t)}), y_{i,j}^{NTP(t)}), \quad (13)$$

where φ_0 is the projection head of the NTP task and CE is the Cross Entropy loss function.

5.2 Contextual Similarity

The purpose of NTP task is to enable the sample to learn the relationships between its internal channels. Meanwhile, we need to ensure the independence between different samples, so we further employ the Contextual Similarity task. The main difference between CS task and Contextual Contrasting task in TS-TCC [12] lies in the design of the augmentation method. To help the framework focus more on the dependencies between channels, we further apply asynchronous permutation strategy to generate negative samples. And we also add the original sample to the self-supervised training to enhance the learning ability.

In this task, we generate several positive samples and negative samples for each sample through augmentation strategy. For positive samples, we adopt the Interval Adjustment Strategy which randomly selects a series of intervals, and then adjust all values by jittering. Further, we also adopt a Synchronous Permutation Strategy which segments the whole time series and disrupts the segment order. It helps to maintain the relations between segments and the generated samples are regarded as positive. For negative samples, we adopt an Asynchronous Permutation Strategy which randomly segments and disrupts the segment order for each channel in different ways. In experiments, for each sample, we use the interval adjustment strategy and the synchronous disorder strategy to generate one positive sample respectively, and we use the asynchronous disorder strategy to generate two negative samples. Assuming that the batch size is B , we can generate extra $4B$ augmented samples. Therefore the total number of the sample in a batch is $5B$.

With the i -th sample x_i in a batch our encoder can obtain its representation $z_i^0 \in \mathbb{R}^{1 \times (C \cdot D)}$. The representations of inputs except itself are $z_i^* \in \mathbb{R}^{(5B-1) \times (C \cdot D)}$, where $z_i^{*,m} \in \mathbb{R}^{1 \times (C \cdot D)}$ is the m -th representation of z_i^* . Among z_i^* , the two positive samples are $z_i^{+,1}, z_i^{+,2} \in \mathbb{R}^{1 \times (C \cdot D)}$. Therefore, for i -th sample, the loss of the CS task can be obtained through the following formula:

$$\ell_{CS} = - \sum_{n=1}^2 \log \frac{\exp(\text{sim}(\varphi_1(z_i^0), \varphi_1(z_i^{+,n}))/\tau)}{\sum_{m=1}^{5B-1} \exp(\text{sim}(\varphi_1(z_i^0), \varphi_1(z_i^{*,m}))/\tau)}, \quad (14)$$

where τ is a hyperparameter, φ_1 is the projection head of the CS task, sim is the cosine similarity.

The final self-supervised loss is the combination of the NTP loss and CS loss as follows:

$$\ell = \alpha_1 \cdot \ell_{NTP} + \alpha_2 \cdot \ell_{CS}, \quad (15)$$

where α_1 and α_2 are hyperparameters.

6 Experiments

6.1 Datasets

To demonstrate the effectiveness of our self-supervised framework, we use the following four datasets from different fields:

- **UCI HAR**³: The UCI HAR dataset [1] contains sensor readings of 6 human activity types. They are collected with a sampling rate of 50 Hz.
- **LSST**⁴: The LSST dataset [2] is an open data to classify simulated astronomical time series data in preparation for observations from the Large Synoptic Survey Telescope.
- **ArabicDigits**⁵: The ArabicDigits dataset [5] contains time series of mel-frequency cepstrum coefficients corresponding to spoken Arabic digits. It includes data from 44 male and 44 female native Arabic speakers.
- **JapaneseVowels**⁵: In the JapaneseVowels dataset [5], several Japanese male speakers are recorded saying the vowels ‘a’ and ‘e’. A ‘12-degree linear prediction analysis’ is applied to the raw recordings to obtain time-series with 12 dimensions.

The detailed information is shown on Table 1.

Table 1. Detailed information of the used datasets.

Dataset	Train	Test	Time	Channel	Class
UCI HAR	7352	2947	128	9	6
LSST	2459	2466	36	6	14
ArabicDigits	6600	2200	93	13	10
JapaneseVowels	270	370	29	12	9

³ <https://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones>

⁴ <http://www.timeseriesclassification.com/description.php?Dataset=LSST>

⁵ <http://www.mustafabaydogan.com/>

6.2 Experimental Settings

For supervised learning, we set $D = 512$, $N = 8$, attention head = 8, batch size = 4, dropout = 0.2. We use Adam optimizer with a learning rate of 1e-4.

For self-supervised learning, we set $K_{NTP} = 10$, $D = 512$, $N = 8$, attention head = 8, batch size = 10, dropout = 0.2, $\tau = 0.2$, $\alpha_1 = 2$, $\alpha_2 = 1$. We use Adam optimizer with a learning rate of 5e-5.

For fine-tuning after self-supervised learning, we train a single fully connected layer on top of the frozen self-supervised pre-trained encoder to evaluate the effect of our self-supervised framework. We set batch size = 4 and use Adam optimizer with a learning rate of 1e-3.

We evaluate the performance using two metrics: Accuracy (ACC) and Macro-F1 score (MF1). Every result is generated by repeating 5 times with 5 different seeds.

We conduct our experiments using PyTorch 1.7 and train models on a NVIDIA GeForce RTX 2080 Ti GPU.

6.3 Baselines

We compare our framework against the following self-supervised methods. It is noted that we apply the default hyperparameters of each compared method from the original paper or the code. The detailed information and the reason why we choose these methods are as followed:

(1) **W2V** [15]: This method employs the idea of word2vec. It combines an encoder based on causal dilated convolutions with a triplet loss and time-based negative sampling. Finally, they train a SVM on top of the frozen self-supervised pre-trained encoder. It achieves great results leveraging unsupervised learning for univariate and multivariate classification datasets. We select $K = 10, 20$ from the original experiments.

(2) **W2V+**: It applies our proposed two tower Transformer-based model as the encoder while training with the pretext task of W2V. For the requirement of the time-wise feature, we replace Channel-aware Transformer with Time-aware Transformer (TaT) which integrates the channel-wise features into the time-wise features in the Aggregate Layer.

(3) **TS-TCC** [12]: This method employs the idea of contrastive learning using a convolutional architecture as encoder. After self-supervised learning, they train a fully connected layer on top of the frozen self-supervised pre-trained encoder. It is the state-of-the-art contrastive learning method in the field of self-supervised of time series.

(4) **TS-TCC+**: It applies TaT as the encoder like W2V+ while training with the pretext task of TS-TCC.

(5) **TST** [34]: This method employs the idea of masked language model using a Transformer-based encoder. In their work, it achieves outstanding results by finetuning the whole encoder after pre-training it. For fair comparison, we freeze the encoder while finetuning.

(6) **TST+**: It applies TaT as the encoder like W2V+ while training with the

pretext task of TST.

(7) **NVP+CS**: To compare with the regression, in our framework we replace Next Trend Prediction with Next Value Predict (NVP) and regard it as a new strong baseline. We select 15% time steps to predict the values of the next time step.

(8) **Supervised**: Supervised learning on both encoder and fully connected layer.

Table 2. Comparison between CaSS and other self-supervised methods. \uparrow mark indicates that the self-supervised result performs better than the supervised result.

Method	HAR		LSST		ArabicDigits		JapaneseVowels	
	ACC	MF1	ACC	MF1	ACC	MF1	ACC	MF1
W2V K=10	90.37±0.34	90.67±0.97	57.24±0.24	36.97±0.49	90.16±0.57	90.22±0.52	97.98±0.40	97.73±0.48
W2V K=20	90.08±0.18	90.10±0.16	53.47±0.75	32.84±1.25	90.55±0.77	90.60±0.75	97.98±0.40	97.89±0.43
W2V+	84.55±0.59	84.34±0.69	54.90±0.58	33.97±0.79	87.59±0.48	87.56±0.49	95.27±0.38	95.25±0.39
TS-TCC	90.74±0.25	90.23±0.29	40.38±0.35	23.93±1.93	95.64±0.37	95.43±0.37	82.25±1.16	82.04±1.17
TS-TCC+	90.87±0.31	90.86±0.30	50.75±0.23	32.87±0.65	96.87±0.34	96.80±0.28	84.36±0.27	84.09±0.21
TST	77.62±2.48	78.05±2.56	32.89±0.04	7.86±1.63	90.73±0.36	90.90±0.33	97.30±0.27	97.47±0.34
TST+	87.39±0.49	87.77±0.11	34.49±0.38	14.62±0.32	96.82±0.23	96.82±0.22	97.87±0.23 \uparrow	97.51±0.22
NVP+CS	92.47±0.20	92.38±0.19	32.42±0.14	6.02±0.30	96.50±0.45	96.51±0.55	96.34±0.43	95.49±0.11
CaSS	92.57±0.24\uparrow	92.40±0.17\uparrow	64.94±0.02	46.11±0.55	97.07±0.20	97.07±0.20	98.11±0.27\uparrow	98.14±0.08\uparrow
Supervised	92.35±0.63	92.40±0.58	66.57±0.38	51.60±1.26	98.07±0.38	98.07±0.38	97.71±0.13	97.56±0.07

6.4 Results and Analysis

Comparison with self-supervised methods The experimental results are shown in Table 2. Overall, our self-supervised learning framework can significantly surpass the previous state-of-the-art methods. Especially in LSST and JapaneseVowels whose time lengths are relative short, methods based on fitting specific values or features like TST, TS-TCC and NVP cannot perform well, while our framework can obtain promising and stable performances. It demonstrates that simple trend predicting is more efficient than regression. W2V and our framework are suitable for both short and long time length datasets, and our performances can significantly surpass W2V. This demonstrates the powerful representation learning ability of our framework. Moreover, our self-supervised framework is shown to be superior to the supervised way in two of four datasets while in the other two datasets can achieve similar results. It also proves that our self-supervised learning framework is capable of learning not only complex characteristics between samples but also within samples.

For a more convincing comparison, we conduct experiments by replacing the origin encoder of the previous methods with our proposed encoder. It helps to offer a more detailed view on the aspects of both encoder and pretext task. It is shown in Table 2 that applying our encoder like TS-TCC+ and TST+ helps to improve their effectiveness, which demonstrates the encoding ability of our encoder. W2V is a method which pays more attention on the local information, so the causal dilated convolutions which only focuses on previous information is

more suitable than W2V+ which encodes the global information. In the aspect of pretext task, when combining NTP task and CS task with our encoder, the self-supervised learning ability is further improved by a large margin compared to other pretext tasks.

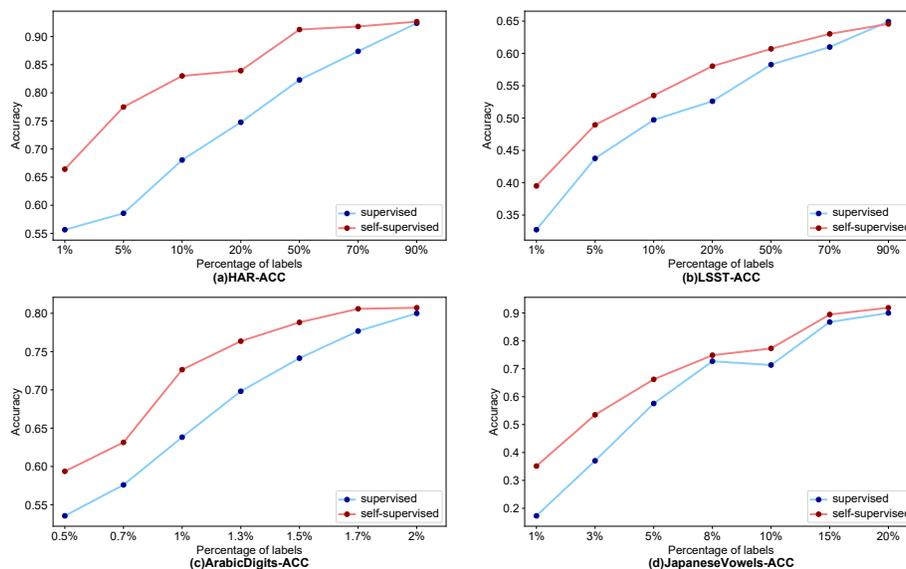


Fig. 4. Few-shot learning results. We report the comparison between our self-supervised framework and the supervised way.

Few-shot learning To further prove the effect of our self-supervised framework, we conduct few-shot learning experiments comparing with the supervised learning. In HAR and LSST datasets, we choose 1%, 5%, 10%, 20%, 50%, 70%, 90% percentage of labeled samples for model training respectively. For ArabicDigits and JapaneseVowels datasets, we adopt a set of smaller percentages in order to compare the performance of few-shot learning more clearly. The results are shown in Figure 4. Among these datasets, our self-supervised framework can significantly surpass the supervised learning by training a single fully connected layer with limited labeled samples.

6.5 Ablation Study

Ablation study on pretext task To analyze the role of each of our pretext tasks, we apply the following variants as comparisons:

- (1) **-NTP**: It only applies Contextual Similarity task to self-supervised learning.
- (2) **-CS**: It only applies Next Trend Prediction task to self-supervised learning.

(3) **-neg augment**: It removes the negative samples generated by the asynchronous disorder strategy.

(4) **reverse neg**: It regards the samples generated by the asynchronous disorder strategy as positive samples.

The experimental results are shown in Table 3. It can be seen that in our pretext tasks, NTP task and CS task can well cooperate with each other. Specifically, NTP task occupies the most important position in small datasets while CS task is more important in datasets with large number of samples. On the one hand, it shows the importance of the internal relationships between channels to multivariate time series. On the other hand, it demonstrates that self-supervised learning should not only focus on the characteristics of the sample itself, but also need to maintain the independence from other samples. In addition, negative sample enhancement can bring a more stable effect to the encoder and further enhance the effect of self-supervised learning method. The conversion of the negative samples into positive samples will lead to the decline of the effect, which shows that in the time series, the time relationships between channels can not be disturbed.

Table 3. Ablation study on pretext task.

Method	HAR		LSST		ArabicDigits		JapaneseVowels	
	ACC	MF1	ACC	MF1	ACC	MF1	ACC	MF1
-NTP	90.96±0.59	90.71±0.59	61.05±1.20	41.30±0.40	95.55±0.36	95.50±0.41	94.87±0.27	94.87±0.22
-CS	76.94±0.66	75.96±0.67	53.00±0.81	33.50±1.97	90.71±0.06	90.72±0.08	97.71±0.13	97.69±0.13
-neg augment	87.06±0.25	86.79±0.23	61.88±0.57	45.21±0.10	96.32±0.32	96.32±0.32	97.98±0.13	97.79±0.16
reverse neg	83.02±0.22	82.65±0.15	63.30±0.65	45.40±0.50	92.62±0.43	92.61±0.44	97.42±0.16	97.22±0.09
CaSS	92.57±0.24	92.40±0.17	64.94±0.02	46.11±0.55	97.07±0.20	97.07±0.20	98.11±0.27	98.14±0.08

Ablation study on encoder To analyze the effect of each component in the encoder, we apply the following variants for comparisons by supervised learning:

(1) **Self Aggregate**: It contains two Transformers with self-attention mechanism to encode time-wise and channel-wise features independently. Finally the time-wise features are fused into channel-wise features through Aggregate Layer.

(2) **Channel Self**: A single Transformer is applied with self-attention mechanism to encode channel-wise features.

(3) **-Aggregate Layer**: The channel-wise features of the last Co-Transformer Layer are applied without fusing the time-wise features.

The experimental results are shown in Table 4. It can be seen that, if time-wise and channel-wise features are only fused in the last aggregate layer without interactions in the previous stage, they cannot be well integrated and bring the loss of information. As contrast, the interactions between the two Transformers can significantly bring the performance improvement. The results of Channel Self illustrate the importance of each channel’s independent time pattern. Finally, the existence of Aggregate Layer can also better integrate the features of time-wise and channel-wise while alleviating the redundancy of features.

Table 4. Ablation study on encoder.

Method	HAR		LSST		ArabicDigits		JapaneseVowels	
	ACC	MF1	ACC	MF1	ACC	MF1	ACC	MF1
Self Aggregate	91.88±0.35	91.85±0.38	66.10±0.24	42.80±0.83	96.64±0.22	96.63±0.23	97.17±0.13	96.94±0.10
Channel Self	93.15±0.91	93.13±0.95	56.45±0.48	31.31±0.50	97.48±0.11	97.47±0.12	97.57±0.27	97.42±0.22
-Aggregate Layer	92.35±0.12	92.36±0.13	63.97±0.59	42.65±0.18	97.68±0.32	97.68±0.32	97.57±0.27	97.38±0.23
CaT	92.35±0.63	92.40±0.58	66.57±0.38	51.60±1.26	98.07±0.38	98.07±0.38	97.71±0.13	97.56±0.07

7 Conclusion

Self-supervised learning is essential for multivariate time series. In this work, we propose a new self-supervised learning framework CaSS to learn the complex representations of MTS. For the encoder, we propose a new Transformer-based encoder Channel-aware Transformer to capture the time-wise and channel-wise features more efficiently. For the pretext task, we propose Next Trend Prediction from the perspective of channel-wise for the first time and combine it with Contextual Similarity task. These novel pretext tasks can well cooperate with our encoder to learn the characteristics. Our self-supervised learning framework demonstrates significant improvement on MTS classification comparing with previous works, and can significantly surpass supervised learning with limited labeled samples.

8 Acknowledgments

This work was supported by the National Key Research and Development Program of China, No.2018YFB1402600.

References

1. Anguita, D., Ghio, A., Oneto, L., Parra Perez, X., Reyes Ortiz, J.L.: A public domain dataset for human activity recognition using smartphones. In: Proceedings of the 21th International European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning. pp. 437–442 (2013)
2. Bagnall, A., Dau, H.A., Lines, J., Flynn, M., Large, J., Bostrom, A., Southam, P., Keogh, E.: The uea multivariate time series classification archive, 2018. arXiv preprint arXiv:1811.00075 (2018)
3. Bai, Y., Wang, L., Tao, Z., Li, S., Fu, Y.: Correlative channel-aware fusion for multi-view time series classification. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 6714–6722 (2021)
4. Banville, H., Albuquerque, I., Hyvärinen, A., Moffat, G., Engemann, D.A., Gramfort, A.: Self-supervised representation learning from electroencephalography signals. In: 2019 IEEE 29th International Workshop on Machine Learning for Signal Processing (MLSP). pp. 1–6. IEEE (2019)
5. Baydogan, M.G.: Multivariate time series classification datasets (2019)
6. Binkowski, M., Marti, G., Donnat, P.: Autoregressive convolutional neural networks for asynchronous time series. In: International Conference on Machine Learning. pp. 580–589. PMLR (2018)

7. Breiman, L.: Random forests. *Machine learning* **45**(1), 5–32 (2001)
8. Che, Z., Purushotham, S., Cho, K., Sontag, D., Liu, Y.: Recurrent neural networks for multivariate time series with missing values. *Scientific reports* **8**(1), 1–12 (2018)
9. Cortes, C., Vapnik, V.: Support-vector networks. *Machine learning* **20**(3), 273–297 (1995)
10. Cui, Z., Chen, W., Chen, Y.: Multi-scale convolutional neural networks for time series classification. *arXiv preprint arXiv:1603.06995* (2016)
11. Deng, H., Runger, G., Tuv, E., Vladimir, M.: A time series forest for classification and feature extraction. *Information Sciences* **239**, 142–153 (2013)
12. Eldele, E., Ragab, M., Chen, Z., Wu, M., Kwok, C.K., Li, X., Guan, C.: Time-series representation learning via temporal and contextual contrasting. In: *International Joint Conference on Artificial Intelligence (IJCAI-21)* (2021)
13. Fan, H., Zhang, F., Gao, Y.: Self-supervised time series representation learning by inter-intra relational reasoning. *arXiv preprint arXiv:2011.13548* (2020)
14. Fawaz, H.I., Forestier, G., Weber, J., Idoumghar, L., Muller, P.A.: Deep learning for time series classification: a review. *Data mining and knowledge discovery* **33**(4), 917–963 (2019)
15. Franceschi, J.Y., Dieuleveut, A., Jaggi, M.: Unsupervised scalable representation learning for multivariate time series. In: *Thirty-third Conference on Neural Information Processing Systems*. vol. 32. Curran Associates, Inc. (2019)
16. Karim, F., Majumdar, S., Darabi, H., Chen, S.: Lstm fully convolutional networks for time series classification. *IEEE access* **6**, 1662–1669 (2017)
17. Kenton, J.D.M.W.C., Toutanova, L.K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of NAACL-HLT*. pp. 4171–4186 (2019)
18. Lin, J., Khade, R., Li, Y.: Rotation-invariant similarity in time series using bag-of-patterns representation. *Journal of Intelligent Information Systems* **39**(2), 287–315 (2012)
19. Lin, L., Xu, B., Wu, W., Richardson, T.W., Bernal, E.A.: Medical time series classification with hierarchical attention-based temporal convolutional networks: A case study of myotonic dystrophy diagnosis. In: *CVPR Workshops*. pp. 83–86 (2019)
20. Liu, C.L., Hsaio, W.H., Tu, Y.C.: Time series classification with multivariate convolutional neural network. *IEEE Transactions on Industrial Electronics* **66**(6), 4788–4797 (2018)
21. Liu, M., Ren, S., Ma, S., Jiao, J., Chen, Y., Wang, Z., Song, W.: Gated transformer networks for multivariate time series classification. *arXiv preprint arXiv:2103.14438* (2021)
22. Ma, F., Chitta, R., Zhou, J., You, Q., Sun, T., Gao, J.: Dipole: Diagnosis prediction in healthcare via attention-based bidirectional recurrent neural networks. In: *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. pp. 1903–1911 (2017)
23. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. In: *1st International Conference on Learning Representations, ICLR* (2013)
24. Oh, J., Wang, J., Wiens, J.: Learning to exploit invariances in clinical time-series data using sequence transformer networks. In: *Machine Learning for Healthcare Conference*. pp. 332–347. PMLR (2018)
25. Poli, M., Park, J., Ilievski, I.: Wattnet: Learning to trade fx via hierarchical spatio-temporal representation of highly multivariate time series. In: *Twenty-Ninth In-*

- ternational Joint Conference on Artificial Intelligence and Seventeenth Pacific Rim International Conference on Artificial Intelligence, IJCAI 2020 (2020)
26. Rasul, K., Sheikh, A.S., Schuster, I., Bergmann, U.M., Vollgraf, R.: Multivariate probabilistic time series forecasting via conditioned normalizing flows. In: International Conference on Learning Representations (2020)
 27. Rußwurm, M., Körner, M.: Self-attention for raw optical satellite time series classification. *ISPRS Journal of Photogrammetry and Remote Sensing* **169**, 421–435 (2020)
 28. Schäfer, P.: The boss is concerned with time series classification in the presence of noise. *Data Mining and Knowledge Discovery* **29**(6), 1505–1530 (2015)
 29. Sun, C., Myers, A., Vondrick, C., Murphy, K., Schmid, C.: Videobert: A joint model for video and language representation learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7464–7473 (2019)
 30. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
 31. Wang, Z., Yan, W., Oates, T.: Time series classification from scratch with deep neural networks: A strong baseline. In: 2017 International joint conference on neural networks (IJCNN). pp. 1578–1585. IEEE (2017)
 32. Xi, X., Keogh, E., Shelton, C., Wei, L., Ratanamahatana, C.A.: Fast time series classification using numerosity reduction. In: Proceedings of the 23rd international conference on Machine learning. pp. 1033–1040 (2006)
 33. Xu, J., Wang, J., Long, M., et al.: Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems* **34** (2021)
 34. Zerveas, G., Jayaraman, S., Patel, D., Bhamidipaty, A., Eickhoff, C.: A transformer-based framework for multivariate time series representation learning. In: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining. pp. 2114–2124 (2021)
 35. Zhang, Y., Wu, Q., Peng, N., Dai, M., Zhang, J., Wang, H., et al.: Memory-gated recurrent networks. In: Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI-21) (2020)
 36. Zhao, B., Lu, H., Chen, S., Liu, J., Wu, D.: Convolutional neural networks for time series classification. *Journal of Systems Engineering and Electronics* **28**(1), 162–169 (2017)
 37. Zheng, Y., Liu, Q., Chen, E., Ge, Y., Zhao, J.L.: Time series classification using multi-channels deep convolutional neural networks. In: International conference on web-age information management. pp. 298–310. Springer (2014)
 38. Zheng, Y., Liu, Q., Chen, E., Ge, Y., Zhao, J.L.: Exploiting multi-channels deep convolutional neural networks for multivariate time series classification. *Frontiers of Computer Science* **10**(1), 96–112 (2016)
 39. Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., Zhang, W.: Informer: Beyond efficient transformer for long sequence time-series forecasting. In: Proceedings of AAAI (2021)