

Improving Position Encoding of Transformers for Multivariate Time Series Classification

Navid Mohammadi Foumani^{1*}, Chang Wei Tan¹, Geoffrey I. Webb¹ and Mahsa Salehi¹

^{1*}Department of Data Science and Artificial Intelligence,
Monash University, Melbourne, VIC, Australia.

*Corresponding author(s). E-mail(s):

navid.foumani@monash.edu;

Contributing authors: chang.tan@monash.edu;

geoff.webb@monash.edu; mahsa.salehi@monash.edu;

Abstract

Transformers have demonstrated outstanding performance in many applications of deep learning. When applied to time series data, transformers require effective position encoding to capture the ordering of the time series data. The efficacy of position encoding in time series analysis is not well-studied and remains controversial, e.g., whether it is better to inject absolute position encoding or relative position encoding, or a combination of them. In order to clarify this, we first review existing absolute and relative position encoding methods when applied in time series classification. We then proposed a new absolute position encoding method dedicated to time series data called time Absolute Position Encoding (tAPE). Our new method incorporates the series length and input embedding dimension in absolute position encoding. Additionally, we propose computationally Efficient implementation of Relative Position Encoding (eRPE) to improve generalisability for time series. We then propose a novel multivariate time series classification (MTSC) model combining tAPE/eRPE and convolution-based input encoding named ConvTran to improve the position and data embedding of time series data. The proposed absolute and relative position encoding methods are simple and efficient. They can be easily integrated into transformer blocks and used for downstream tasks such as forecasting, extrinsic regression, and anomaly detection. Extensive experiments on 32 multivariate time-series datasets

show that our model is significantly more accurate than state-of-the-art convolution and transformer-based models. Code and models are open-sourced at <https://github.com/Navidfoumani/ConvTran>.

Keywords: Multivariate Time Series Classification, Transformers, Position Encoding

1 Introduction

A time series is a time-dependent quantity recorded over time. Time series data can be univariate, where only a sequence of values for one variable is collected; or multivariate, where data are collected on multiple variables. There are many applications that require time series analysis, such as human activity recognition [1], diagnosis based on electrocardiogram (ECG), electroencephalogram (EEG), and systems monitoring problems [2]. Many of these applications are inherently multivariate in nature – various sensors are used to measure human’s activities; EEGs use a set of electrodes (channels) to measure brain signals at different locations of the brain. Hence, multivariate time-series analysis methods such as classification and segmentation are of great current interest [3–5].

Convolutional neural networks (CNNs) have been widely employed in time series classification [4, 5]. Many studies have shown that convolution layers tend to have strong generalization with fast convergence due to their strong inductive bias [6]. While CNN-based models are excellent for capturing local temporal/spatial correlations, these models cannot effectively capture and utilize long-range dependencies. Also, they only consider the local order of data points in a time series rather than the order of all data points globally. Due to this, many recent studies have used recurrent neural networks (RNN) such as LSTMs to capture this information [7]. However, RNN-based models are computationally expensive, and their capability in capturing long-range dependencies are limited [8, 9].

On the other hand, *attention models* can capture long-range dependencies, and their broader receptive fields provide more contextual information, which can improve the models’ learning capacity. Not surprisingly, with the success of attention models in natural language processing [8, 10], many previous studies have attempted to bring the power of attention models into other domains such as computer vision [11] and time series analysis [9, 12, 13].

The transformer’s core is self-attention [8], which is capable of modeling the relationship of input time series. Self-attention, however, has a limitation — it cannot capture the ordering of input series. Hence, adding explicit representations of position information is especially important for the attention since the model is otherwise entirely invariant to input order, which is undesirable for modeling sequential data. This limitation is even worse in time series

data since, unlike image and text, which use Word2Vec-like embedding, time series data has less informative data context.

There are two main methods for encoding positional information in transformers: absolute and relative. Absolute methods, such as those used in [8, 10], assign a unique encoding vector to each position in the input sequence based on its absolute position in the sequence. These encoding vectors are combined with the input encoding to provide positional information to the model. On the other hand, relative methods [14, 15] encode the relative distance between two elements in the sequence, rather than their absolute positions. The model learns to compute the relative distances between any two positions during training and looks up the corresponding embedding vectors in a pre-defined table to obtain the relative position embeddings. These embeddings are used to directly modify the attention matrix. Position encoding has been verified to be effective in natural language processing and computer vision [16]. However, in time series classification, the efficacy is still unclear.

The original absolute position encoding is proposed for language modeling, where high embedding dimensions like 512 or 1024 are usually used for position embedding of input with a length of 512 [8]. But, for time series tasks, embedding dimensions are relatively low, and the series might have a variety of lengths (ranging from very low to very high). In this paper, for the first time, we study the efficiency (i.e. how well resources are utilized) and the effectiveness (i.e. how well the encodings achieve their intended purpose) of existing absolute and relative position encodings for time series data. We then show that the existing absolute position encodings are ineffective with time series data. We introduce a novel time series-specific absolute position encoding method that takes into account the series embedding dimension and length. We show that our new absolute position encoding outperforms the existing absolute position encodings in time series classification tasks.

Additionally, since the existing relative position encodings have large memory overhead and they require a large number of parameters to be trained, in time series data it is very likely they overfit. We propose a novel computationally efficient implementation of relative position encoding to improve their generalisability for time series. We show that our new relative position encoding outperforms the existing relative position encodings in time series classification tasks. We then propose a novel time series classification model based on the combination of our proposed absolute/relative position encodings named ConvTran to improve the position embedding of time series data. We further enriched the data embedding of time series using CNN rather than linear encoding. Our extensive experiments on 32 benchmark datasets show ConvTran is significantly more accurate than the previous state-of-the-art in deep learning models for time series classification (TSC). We believe our novel position encodings can boost the performance of other transformer-based TSC models.

2 Related Work

In this section, we briefly discuss the state-of-the-art multivariate time series classification (MTSC) algorithms, as well as CNN and attention-based models that have been applied to MTSC tasks. We refer interested readers to the corresponding papers or the recent survey on deep learning for time series classification [17] for a more detailed description of these algorithms and models.

2.1 State-of-the-art MTSC Algorithms

Many MTSC algorithms have been proposed in recent years [2, 4, 5], where many of them are adapted from their univariate version. A recent survey [5] evaluated most of the existing MTSC algorithms on the UEA MTS archive, that consists of 26 equal-length time series datasets. This benchmark includes a few deep learning as well as non-deep learning approaches. This survey concluded that there are four main state of the art methods. These are ROCKET [18], HIVE-COTE [19], CIF [20] and Inception-Time [21].

ROCKET [18] is a scalable TSC algorithm that uses 10,000 random convolution kernels to extract 2 features from each input time series, creating 20,000 features for each time series. Then a linear model is used for classification, such as ridge or logistic regression. Mini-ROCKET [22] is an extension of ROCKET with some slight modifications to the feature extraction process. It is significantly more scalable than ROCKET and uses only 10,000 features without compromising accuracy. Multi-ROCKET [23] extends Mini-ROCKET by leveraging the first derivative of the series as well as extracting 4 features per kernel. It is significantly more accurate than both ROCKET and Mini-ROCKET on 128 univariate TSC tasks. Note that neither Mini-ROCKET nor Multi-ROCKET has previously been benchmarked on the UEA MTS archive. The adaptation for multivariate time series for ROCKET, Mini-ROCKET and Multi-ROCKET is done by randomly selecting different channels of the time series for each convolutional kernel.

The Canonical Interval Forest (CIF) [20] is an interval based classifier. It first extracts 25 features from random intervals of the time series and builds a time series forest with 500 trees. It is an algorithm initially designed for univariate TSC and was adapted to multivariate TSC by expanding the random interval search space, where an interval is defined as a random dimension of the time series.

The Hierarchical Vote Collective of Transformation-based Ensembles (HIVE-COTE) is a meta ensemble for TSC. It forms its ensemble from classifiers of multiple domains. Since its introduction in 2016, HIVE-COTE has gone through a few iterations. The version used in the MTSC benchmark [5] comprised of 4 ensemble members – Shapelet Transform Classifier (STC), Time Series Forest (TSF), Contractable Bag of Symbolic Fourier Approximation Symbols (CBOSS) and Random Interval Spectral Ensemble (RISE), each

of them being the state of the art in their respective domains. Since these algorithms were designed for univariate time series, the adaptation for multivariate time series is not easy. Hence, they were adapted for multivariate time series through ensembling over all the models built on each dimension independently. This means that they are computationally very expensive especially when the number of channels is large. Recently, the latest HIVE-COTE version, HIVE-COTEv2.0 (HC2) was proposed [24]. It is currently the most accurate classifier for both univariate and multivariate TSC tasks [24]. Despite being the most accurate on 26 benchmark MTSC datasets, that are relatively small, HC2 is not scalable to either large datasets with long time series or datasets with many channels.

2.2 CNN Based Models

CNNs are popular deep learning architectures for MTSC due to their ability to extract latent features from the time series data efficiently. Fully Convolutional Neural Network (FCN) and Residual Network (ResNet) were proposed in [25] and evaluated in [4]. FCN is a simple convolutional network that does not contain any pooling layers in convolution blocks. The output from the last convolution block is averaged with a Global Average Pooling (GAP) layer and passed to a final softmax classifier. ResNet is one of the deepest architectures for MTSC (and TSC in general), containing three residual blocks followed by a GAP layer and a softmax classifier. It uses residual connections between blocks to reduce the vanishing gradient effect in deep learning models. ResNet was one of the most accurate deep learning TSC architectures on 85 univariate TSC datasets [3, 4]. It was also proven to be an accurate deep learning model for MTSC [4, 5].

Inception-Time is the current state-of-the-art deep learning model for both univariate TSC and MTSC [5, 21]. Inception-Time is an ensemble of five randomly initialised inception network models that each consists of two blocks of inception modules. Each inception module first reduces the dimensionality of a multivariate time series using a bottleneck layer with length and stride of 1 while maintaining the same length. Then, 1D convolutions of different lengths are applied to the output of the bottleneck layer to extract patterns at different sizes. A max pooling layer followed by a bottleneck layer are also applied to the original time series to increase the robustness of the model to small perturbations. Residual connections are also used between each inception block to reduce the vanishing gradient effect. The output of the second inception block is passed to a GAP layer before feeding into a softmax classifier.

Recently, Disjoint-CNN [26] shows that factorization of 1D convolution kernels into disjoint temporal and spatial components yields accuracy improvements with almost no additional computational cost. Applying disjoint temporal convolution and then spatial convolution behaves similarly to the “Inverted Bottleneck” [27]. Like the Inverted Bottleneck, the temporal convolutions expand the number of input channels, and spatial convolutions later

project the expanded hidden state back to the original size to capture the temporal and spatial interaction.

2.3 Attention Based Models

Self-attention has been demonstrated to be effective in various natural language processing tasks due to its higher capacity and superior ability to capture long-term dependencies in text [8]. Recently, it has also been shown to be effective for time series classification tasks. Cross Attention Stabilized Fully Convolutional Neural Network (CA-SFCN) [9] has applied the self-attention mechanism to leverage the long-term dependencies for the MTSC task. CA-SFCN combines FCN and two types of self-attention - temporal attention (TA) and variable attention (VA), which interact to capture both long-range temporal dependencies and interactions between variables. With evidence that multi-headed attention dominates self-attention, many models try to adapt it to the MTSC domain. Gated Transformer Networks (GTN) [28], similar to CA-SFCN, use two-tower multi-headed attention to capture discriminative information from the input series. They merge the output of two towers using a learnable matrix named gating.

Inspired by the development of transformer-based self-supervised learning like BERT [13], many models try to adopt the same structure for time series classification [12, 13]. BERT-inspired Neural Data Representations (BENDER) replace the word2vec encoder in BERT with the wav2vec to leverage the same structure for time series data. BENDER shows that if we have a massive amount of EEG data, the pre-trained model can be used effectively to model EEG sequences recorded with differing hardware. Similarly, Voice-to-Series with Transformer-based Attention (V2Sa) uses a large-scale pre-trained speech processing model for downstream problems like time series classification problems [29]. Recently, a Transformer-based Framework (TST) was also introduced to adopt vanilla transformers to the multivariate time series domain [12]. TST uses only the encoder part of transformers and pre-train it with proportionally masked data in an unsupervised manner.

3 Background

This section provides a basic definition of self-attention and an overview of current position encoding models. Note that position encoding refers to the method that integrates position information, e.g., absolute or relative. Position embedding refers to a numerical vector associated with position encoding.

3.1 Problem Description and Notation

Given a time series dataset X with n samples, $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where $\mathbf{x}_t = \{x_1, x_2, \dots, x_L\}$ is a d_x -dimensional time series and L is the length of time series, $\mathbf{x}_t \in \mathbb{R}^{L \times d_x}$, and the set of relevant response labels $Y = \{y_1, y_2, \dots, y_n\}$,

$y_t \in \{1, \dots, c\}$ and c is the number of classes. The aim is to train a neural network classifier to map set X to Y .

3.2 Self-Attention

The first attention mechanisms were proposed in the context of natural language processing [30]. While they still relied on a recurrent neural network at its core, Vaswani et al. [8] proposed a transformer model that relies on attention only. Transformers map a query and a set of key-value pairs to an output. More specifically, for an input series, $\mathbf{x}_t = \{x_1, x_2, \dots, x_L\}$, self-attention computes an output series $\mathbf{z}_t = \{z_1, z_2, \dots, z_L\}$ where $z_i \in \mathbb{R}^{d_z}$ and is computed as a weighted sum of input elements:

$$z_i = \sum_{j=1}^L \alpha_{i,j} (x_j W^V) \quad (1)$$

Each coefficient weight $\alpha_{i,j}$ is calculated using softmax function:

$$\alpha_{i,j} = \frac{\exp(e_{ij})}{\sum_{k=1}^L \exp(e_{ik})} \quad (2)$$

where e_{ij} is an attention weight from positions j to i and is computed using a scaled dot-product:

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K)^T}{\sqrt{d_z}} \quad (3)$$

The projections $W^Q, W^K, W^V \in \mathbb{R}^{d_x \times d_z}$ are parameter matrices and are unique per layer. Instead of computing self-attention once, Multi-Head Attention (MHA) [8] does so multiple times in parallel, i.e., employing h attention heads. A linear transformation is applied to the attention head outputs and concatenated into the standard dimensions.

3.3 Position Encoding

The self-attention layer cannot preserve time series positional information in the transformer architecture since the transformer contains no recurrence and convolution. However, the local positional information, i.e., the ordering of time series, is essential. The practical approach in transformer-based methods involves using multiple encoding [16, 31, 32], such as absolute or relative positional encoding, to enhance the temporal context of time-series inputs.

3.3.1 Absolute Position Encoding

The original self-attention considers the absolute position [8], and adds the absolute positional embedding $P = (p_1, \dots, p_L)$ to the input embedding x as:

$$x_i = x_i + p_i \quad (4)$$

where the position embedding $p_i \in \mathbb{R}^{d_{model}}$. There are several options for absolute positional encodings, including the fixed encodings by sine and cosine functions with different frequencies called *VanillaAPE* and the learnable encodings through trainable parameters (we refer it as *Learn* method) [8, 10].

By using sine and cosine for fixed position encoding, the d_{model} -dimensional embeddings of i_{th} time step position can be represented by the following equation:

$$p_i(2k) = \sin i\omega_k \quad p_i(2k+1) = \cos i\omega_k \quad \omega_k = 10000^{-2k/d_{model}} \quad (5)$$

where k is in the range of $[0, \frac{d_{model}}{2}]$, d_{model} is the embedding dimension and ω_k is the frequency term. Variations in ω_k ensure that no positions $< 10^4$ are assigned similar embeddings.

3.3.2 Relative Position Encoding

In addition to the absolute position embedding, recent studies in natural language processing and computer vision also consider the pairwise relationships between input elements, i.e., relative position [14, 15]. This type of method encodes the relative distance between the input elements x_i and x_j into vectors $p_{i,j}^Q, p_{i,j}^K, p_{i,j}^V \in \mathbb{R}^{d_z}$. The encoding vectors are embedded into the self-attention module, which modifies Equation 1 and Equation 3 as

$$z_i = \sum_{j=1}^L \alpha_{i,j} (x_j W^V + p_{i,j}^V) \quad (6)$$

$$e_{ij} = \frac{(x_i W^Q + p_{i,j}^Q)(x_j W^K + p_{i,j}^K)^T}{\sqrt{d_z}} \quad (7)$$

By doing so, the pairwise positional relation is trained during transformer training.

Shaw et al. [14] proposed the first relative position encoding for self-attention. Relative positional information is supplied to the model on two levels: values and keys. First, relative positional information is included in the model as an additional component to the keys. The softmax operation Equation 3 remains unchanged from vanilla self-attention. Lastly, relative positional information is resupplied as a sub-component of the values matrix. Besides, the authors believe that relative position information is not useful beyond a certain distance, so they introduced a clip function to reduce the number of parameters. Encoding is formulated as follows to consider the distance between inputs i and j in computing their attention:

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K + p_{clip(i-j,k)}^K)^T}{\sqrt{d_z}} \quad (8)$$

$$z_i = \sum_{j=1}^L \alpha_{i,j} (x_j W^V + p_{clip(i-j,k)}^V) \quad (9)$$

$$\text{clip}(x, k) = \max(-k, \min(k, x)) \quad (10)$$

Where p^V and p^K are the trainable weights of relative position encoding on values and keys, respectively. $P^V = (p_{-k}^V, \dots, p_k^V)$ and $P^K = (p_{-k}^K, \dots, p_k^K)$ where $p_i^V, p_i^K \in \mathbb{R}^{d_z}$. The scalar k is the maximum relative distance.

However, this technique (Shaw) is not memory efficient. As can be seen in Equation 8, it requires $O(L^2d)$ memory due to the additional relative position encoding. Huang et al. [15] introduced a new method (in this paper it is called *Vector* method) of computing relative positional encoding that reduces its intermediate memory requirement from $O(L^2d)$ to $O(Ld)$ using skewing operation [15]. According to this paper, the authors also dropped the additional relative positional embedding corresponding to the value term and focused only on the key component. Encoding is formulated as follows:

$$e_{ij} = \frac{(x_i W^Q)(x_j W^K)^T + S^{rel}}{\sqrt{d_z}} \quad (11)$$

$$S^{rel} = \text{Skew}(W^Q P) \quad (12)$$

Where *Skew* procedure use padding, reshaping and slicing to reduce the memory requirement [15]. In Table 1 we provided a summary of the parameter sizes, memory, and computation complexities of various position encoding methods (including our proposed ones in this paper) for comparison purposes.

4 Position Encoding of Transformers for MTSC

We design our position encoding methods to examine several aspects which are not well studied in prior transformers-based time series classification work (see the analysis in Sec 5.4).

As a first step, we propose a new absolute position encoding method dedicated to time series data called **time Absolute Position Encoding** (tAPE). tAPE incorporates the series length and input embedding dimension in absolute position encoding. We then introduce **efficient Relative Position Embedding** (eRPE) to explore the independent encoding of positions from the input encodings. After that, to study the integration of eRPE into a transformer model, we compare different integration of position information to the attention matrix; finally, we provide an efficient implementation for our methods.

4.1 Time Absolute Position Encoding (tAPE)

Absolute position encoding was originally proposed for language modeling, where high embedding dimensions like 512 or 1024 are usually used for position embedding of input with a length of 512 [8]. Fig.1a shows the dot product between two sinusoidal positional embedding whose distance is K using Equation 5 with various embedding dimensions. Clearly, higher embedding dimensions, such as 512 (red thick line), can better reflect the similarity between various positions. As shown in Fig.1a using 64 or 128 as embedding

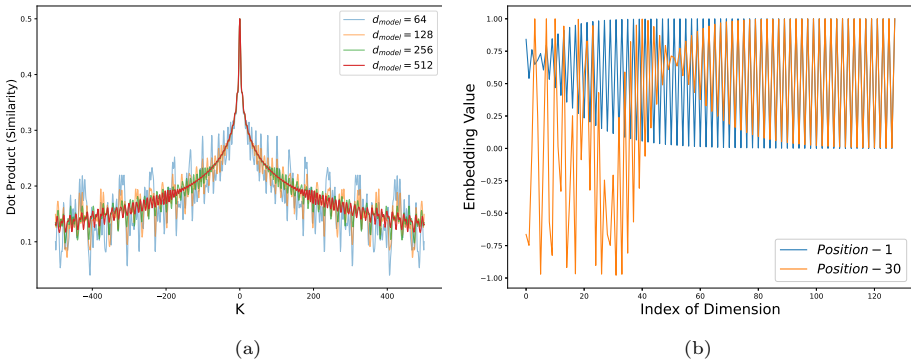


Fig. 1 Sinusoidal absolute position encoding. a) The dot product of two sinusoidal position embeddings whose distance is K with various embedding dimensions. b) 128 dimension sinusoidal positional encoding curves for positions 1 and 30 in a series of length 30.

dimensions (thin blue and orange lines, respectively), the dot product does not always decrease as the distance between two positions increases. We call this the *distance awareness property*, which disappears when lower embedding dimensions, such as 64, are used for position encoding.

While high embedding dimensions show a desirable monotonous decrease trend when the distance between two positions increases (see red line in Fig. 1a), they are not suitable for encoding time series datasets. The reason is that most time series datasets have relatively low input dimensionality (e.g., 28 out of 32 datasets have less than 64 input dimension), and higher embedding dimensions may yield inferior model throughput due to extra parameters (increasing the chances of overfitting the model).

On the other hand, in low embedding dimensions, the similarity value between two random embedding vectors is high, making the embedding vectors very similar to each other. In other words, we cannot fully utilise the embedding vector space to differentiate between two positions. Fig. 1b depicts the embedding vectors of the first and last position embedding for the embedding dimension equals 128 and length equals 30. In this figure, almost half of the embedding vectors are the same. This is called the *anisotropic phenomenon* [33]. The anisotropic phenomenon makes the position encoding to be ineffective in low embedding dimensions as embedding vectors become similar to each other as it is shown in Fig. 1a (the blue line).

Hence, we require a position embedding for time series that has distance awareness while simultaneously being isotropic. In order to incorporate distance awareness, we propose to use the time series length in Equation 5. In this equation, ω_k refers to the frequency of the sine and cosine functions from which the embedding vectors are generated. Without our modification, as series length L increases the dot product of positions becomes ever less regular,

resulting in a loss of distance awareness. By incorporating the length parameter in the frequency terms in both sine and cosine functions in Equation 5, the dot product remains smoother with a monotonous trend.

As the embedding dimension d_{model} value increases, it is more likely the vector embeddings are sampled from low-frequency sinusoidal functions, which results in the anisotropic phenomenon. To alleviate this, we incorporate the d_{model} parameter into the frequency term in both sine and cosine functions in Equation 5. We propose a novel absolute position encoding for time series called tAPE in which ω_k^{new} takes into account the input embedding dimension and length as follows:

$$\begin{aligned}\omega_k &= 10000^{-2k/d_{model}} \\ \omega_k^{new} &= \frac{\omega_k \times d_{model}}{L}\end{aligned}\quad (13)$$

where L is the series length and d_{model} shows the embedding dimension.

Our new tAPE position encoding is compared with a vanilla sinusoidal position encoding to provide further illustration. Using $d_{model} = 128$ dimension vector, Figs 2a-b show the dot product (similarity) of two positions with a distance of K for series with of length $L = 1000$ and $L = 30$ respectively. As depicted in Fig 2a, in vanilla APE, only the closest positions in the series have a monotonous decreasing trend, and approximately from a distance 50 onwards ($|K| > 50$) on both sides, the decreasing similarity trend becomes less apparent as the distance between two positions in the time series increases. However, tAPE has a more stable decreasing trend and more steadily reflects the distance between two positions. Meanwhile, Fig 2b shows the embedding vectors of tAPE are less similar to each other compared to vanilla APE. This is due to better utilising the embedding vector space to differentiate between two positions as we discussed earlier.

Note in Equation 13 our ω_k^{new} will obviously be equal to the ω_k in vanilla APE when $d_{model} = L$ and the encodings of tAPE and vanilla APE will be the same. However, if $d_{model} \neq L$, tAPE will encode the positions in series more effectively than vanilla APE due to the two properties we discussed earlier. Fig 2a shows a case in which $d_{model} < L$ and Fig 2b shows a case in which $d_{model} > L$ and in both cases tAPE utilises embedding space to provide an isotropic encoding, while holding the distance awareness property. In other words, tAPE provides a balance between these two properties in its encodings. The superiority of tAPE compared to vanilla APE and learned APE on various length time series datasets is shown in the experimental results section.

4.2 Efficient Relative Position Encoding (eRPE)

There are multiple extensions of the abovementioned Section 3.3.2 relative position embeddings in machine translation and computer vision [16, 31, 32]. However, input embeddings are the basis for all previous methods of relative position encoding (adding or multiplying the position matrices to the query,

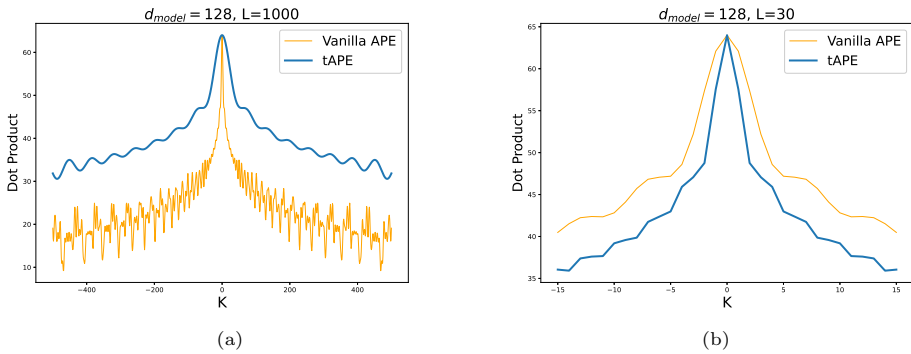


Fig. 2 Comparing dot product between two positions whose distance is K in a time series using tAPE and vanilla APE with $d_x = 128$ dimension vector for series of length a) $L = 1000$ b) $L = 30$.

key, and value matrices). In this study, we introduce an efficient model of relative position encoding independent of input embeddings.

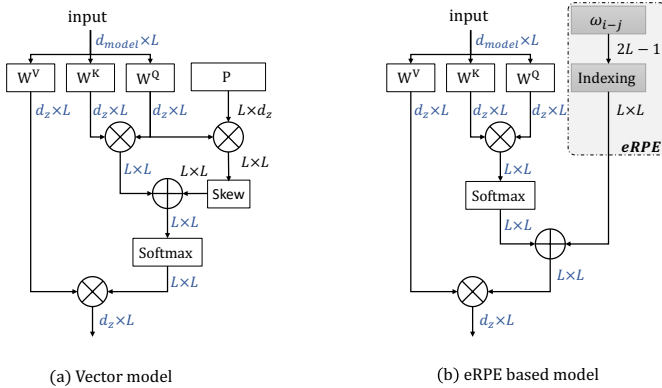


Fig. 3 Self-attention modules with relative position encoding using scalar and vector parameters. Newly added parts are depicted in grey.

In particular, we propose the following formulation:

$$\alpha_i = \sum_{j \in L} \left(\frac{\exp(e_{i,j})}{\underbrace{\sum_{k \in L} \exp(e_{i,k})}_{A_{i,j}}} + w_{i-j} \right) x_j \quad (14)$$

where L is series length, $A_{i,j}$ is attention weight and w_{i-j} is a learnable scalar (i.e., $w \in \mathbb{R}^{O(L)}$) and represent the relative position weight between positions i and j .

It is worth comparing the strengths and weaknesses of relative position encodings and attention to determine what properties are more desirable for relative position encoding of time series data. Firstly, the relative position embedding w_{i-j} is an input-independent parameter with static values, whereas an attention weight $A_{i,j}$ is dynamically determined by the representation of the input series. In other words, attention adapts to input series via a weighting strategy (input-adaptive weighting [8]). Input-adaptive-weighting enables models to capture the complicated relationships between different time points, a property that we desire most when we want to extract high-level concepts in time series. This can be for instance the seasonality component in time series. However, when we have limited size data we are at a greater risk of overfitting when using attention.

Secondly, relative position embedding w_{i-j} takes into account the relative shift between positions i and j and not their values. This is similar to translation equivalence property of convolution, which has been shown to enhance generalization [6]. We propose to consider the notation of w_{i-j} as a scalar rather than a vector to enable the translation equivalency without blowing up the number of parameters. In addition, the scalar representation of w provides the benefit that the value of w_{i-j} for all (i, j) can be subsumed within the pairwise dot-product attention function, resulting in minimal additional computation (see subsection 4.2.1). We call our proposed efficient relative position encoding as eRPE.

Theoretically, there are many possibilities for integrating relative position information into the attention matrix, but we empirically found that attention models perform better when we add the relative position to the model after applying the softmax to the attention matrix as shown in Equation 14. We presume this is because the position values will be sharper without the softmax. And sharper position embeddings seems to be beneficial in TSC task as it emphasizes more on informative relative positions for classification compared to existing models in which softmax is applied to relative position embeddings.

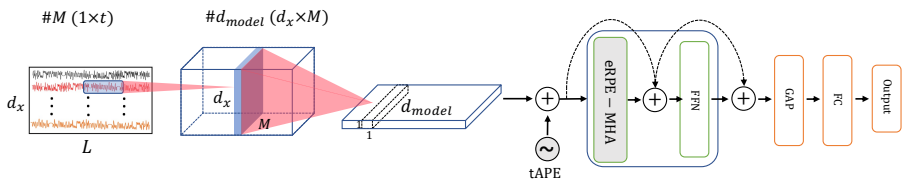
4.2.1 Efficient Implementation: Indexing

To implement the efficient version of eRFE in Equation 14 for input time series with a length of L , for each head, we create a trainable parameter w of size $2L - 1$, as the maximum distance is $2L - 1$. Then for two position indices i and j , the corresponding relative scalar is w_{i-j+L} where indexes start from 1 instead of 0 (1-base index). Accordingly, we need to index L^2 elements from $2L - 1$ vector.

On GPU, a more efficient way to index is to use `gather`, which only requires memory access. At inference time, indexing the L^2 elements from $2L - 1$ vector can be pre-computed and cached to increase the processing speed further. As shown in Table 1, our proposed eRPE is more efficient in terms of both memory

Table 1 Comparing the parameter sizes, memory, and computation complexities of various position encoding methods. In our implementation d_z is equal to d_{model}

	Method	Parameter	Memory	Complexity
	tAPE	None	Ld_{model}	Ld_{model}
Absolute	Vanilla APE [8]	None	Ld_{model}	Ld_{model}
	Learn [10]	Ld_{model}	Ld_{model}	Ld_{model}
	Shaw [14]	$(2L - 1)d_z$	$L^2d_z + L^2$	L^2d_z
Relative	Vector [15]	Ld_z	$Ld_z + L^2$	L^2d_z
	eRPE	$2L - 1$	$L + L^2$	L^2

**Fig. 4** Overall Architecture of the ConvTran Model

and time complexities compared to the existing relative position encoding methods in the literature.

4.3 ConvTran

Now we look at how we can utilize our new position encodings method to build a time series classification network. According to the earlier discussion, global attention has a quadratic complexity w.r.t. the series length. This means that if we directly apply the proposed attention in Equation 14 to the raw time series, the computation will be excessively slow for long time series. Hence, we first use convolutions to reduce the series length and then apply our proposed position encodings once the feature map has been reduced to a less computationally intense size. See Fig. 4 where convolution blocks comes as a first component preceded by attention blocks.

Another benefit of using convolutions is that convolutions operations are very well-suited to capture local patterns. By using convolutions as the first component in our architecture we can capture any discriminative local information that exists in raw time series.

As Shown in Fig. 4, as the first step in the convolution layers, M temporal filters are applied to the input data. In this step, the model extracts temporal patterns in the input series. Next, the output of temporal filtering is convolved with d_{model} spatial $d_x \times M$ shape filters to capture the correlations between variables in multivariate time series and construct d_{model} size input embeddings. Such disjoint temporal and spatial convolution is similar to “Inverted Bottleneck” in [27]. It first expands the number of input channels and then squeezes them. A key reason for this choice is that the Feed Forward

Network (FFN) in transformers [8] also expands on the input size and later projects the expanded hidden state back to the original size to capture the spatial interactions.

Before feeding the input embedding to the transformer block, we add the tAPE-generated position embedding to the input embedding vector so that the model can capture the temporal order of the time series. The size of the embedding vector is d_{model} , which is the same as the input embedding. Inside the multi-head attention, the inputs with the $L \times d_{model}$ dimension are first converted to $L \times d_z \times 3$ shape using a linear layer to get the **qkv** matrix in which d_z indicates the model dimension and defined by the user. Each of the three matrices of shape $L \times d_z$ represents the Query (**q**), Key (**k**) and Value (**v**) matrices. These **q**, **k**, and **v** matrices are reshaped to $h \times L \times d_z/h$ to represent the h attention heads. Each of these attention heads can be responsible for capturing different patterns in time series. For instance, one attention head can attend to the non-noisy data, another head can attend to the seasonal component and another to the trend. Once we have the q , k , and v matrices, we finally perform the attention operation inside the Multi-Head attention block using Equation 14.

According to Equation 14 the eRPE with the same shape of $L \times L$ is also added to the attention output. We consider the notation of w_{i-j} as a scalar (i.e., $w \in R^{O(L)}$) to enable the global convolution kernel without increasing the number of parameters. The relative position embedding enables the model to learn not only the order of time points, but also the relative position of pairs of time points, which can capture richer information than other position embedding strategies.

The FFN, is a multi-layer perceptron block consisting of two linear layers and Gaussian Error Linear Units (GELUs) as an activation function. The outputs from the FFN block are again added to the inputs (via skip connection) to get the final output from the transformer block. Finally, just before the fully connected layer, max-pooling and global average pooling (GAP) are applied to the output of the last layer's ELU activation function, which gives a more translation-equivalence model.

5 Experimental Results

In this section, we evaluate the performance of our ConvTran model on the UEA time series repository [2] and two large multivariate time series datasets and compare it with the state-of-the-art models. All of our experiments were conducted using the PyTorch framework in Python on a computing system consisting of a single Nvidia A5000 GPU with 24GB of memory and an Intel(R) Core(TM) i9-10900K CPU. To promote reproducibility, we have provided our source code and more experimental results online ¹.

We have divided our experiments into four parts. First, we present an ablation study on various position encodings. Then, we demonstrate that our

¹<https://github.com/Navidfoumani/ConvTran>

ConvTran model outperforms existing CNN and transformer-based models. Next, we compare the performance of ConvTran with four state-of-the-art MTSC algorithms (including both deep learning and non-deep learning categories) identified in [5, 24]. We report the results provided on the archive website² for HiveCote2, CIF, ROCKET, and Inception-Time on 26 out of 30 UEA datasets only in Section 5.6. Finally, we evaluate the efficiency and effectiveness of ConvTran by comparing it with the current state-of-the-art model, ROCKET.

5.1 Datasets

UEA Repository The archive consists of 30 real-world multivariate time series data from a wide range of applications such as Human Activity Recognition, Motion classification, and ECG/EEG classification [2]. The number of dimensions ranges from two dimensions to 1345 dimensions. The length of the time series ranges from 8 to 17,984. The datasets also have a train size ranging from 12 to 25000.

Ford Challenge This dataset is obtained from the Kaggle challenge website³. It includes measurements from total of 600 real-time driving sessions where each driving session takes 2 minutes and sampled with 100ms rate. Also, the trials are samples from 100 drivers of both genders, and of different ages. The training data file consists of 604,329 data points each belongs to one of 500 trials. The test file contains 120,840 data points belonging to 100 trials. While each data point comes with a label in 0,1 and also contains 8 physiological, 12 environmental, and 10 vehicular features that are acquired while driving.

Actitracker human Activity Recognition This dataset describes six daily activities which are collected in a controlled laboratory environment. The activities include “Walking”, “Jogging”, “Stairs”, “Sitting”, “Standing”, and “Lying Down” which are recorded from 36 users collected using a cell phone in their pocket. Data has 2,980,765 samples with 3 dimensions, subject-wise split into train and test sets, and a sampling rate of 20Hz [1].

5.2 Evaluation Procedure

We use the classification accuracy as the overall metric to compare different models. Then we rank each model based on its classification accuracy per dataset. The most accurate model is assigned a rank of 1 and the worse performing model is assigned the highest rank. The average ranking is taken in case of ties. Then the average rank for each model is computed across all datasets in the repository.

This gives a direct general assessment of all the models: the lowest rank corresponds to the method that is the most accurate on average. The average ranking for each model is presented in the form of critical difference diagram

²<https://timeseriesclassification.com/HC2.php>

³<https://www.kaggle.com/c/stayalert>

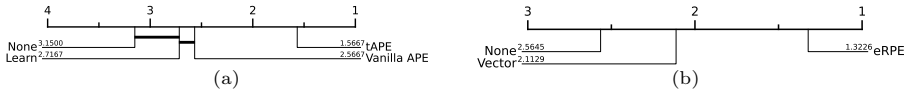


Fig. 5 Critical difference diagram of various position encoding over thirty datasets for the UEA MTSC archive based on average accuracies: a) Various absolute position encodings, b) Various relative position encodings. The lowest rank corresponds to the method that is the most accurate on average.

[34], where models in the same clique (the black bar in the diagram) are not statistically significant. For the statistical test, we used the Wilcoxon signed-rank test with Holm correction as the post hoc test to the Friedman test [34].

5.3 Parameter Setting

Adam optimization is used simultaneously with an early stopping method based on validation loss. We use the default setting for other models. We set the default value for the number of temporal and spatial filters to 64 and set the length of the temporal filters to 8. The width of the spatial convolutions are set equal to the input dimensions [26].

Similar to TST, the transformers based model for MTSC [12], and default transformers block [8], we use 8 heads to capture the varieties of attention from input series. The dimension of transformers encoding is set to $d_{model} = d_z = 64$ and FFN in transformers block expands the input size by 4x and later projects the 4x-wide hidden state back to the original size.

5.4 Ablation Study on Position Encoding

In this section, firstly we compare our proposed tAPE with the existing absolute position encodings. Secondly, we compare our proposed eRPE with the existing relative position encoding methods. As a final step, we combined tAPE and eRPE into a single framework and compare it with all possible combinations of absolute and relative position encodings.

For this ablation study we run a single-layer transformer five times on all 30 UEA benchmark datasets for classification. Fig.5a illustrates the critical difference diagram of a single-layer transformer with different absolute position encodings. Note in critical difference diagram methods grouped by a black line are not significantly different from each other. In Fig.5, *None* is the model without any position encoding, *Learn* is the model with learning absolute position encoding parameters [10], *Vanilla APE* is the vanilla sinusoidal function-based encoding [8], *Vector* is the vector-based implementation of input-dependent relative position embedding [15], and our proposed models showed as *tAPE* and *eRPE*.

As depicted in Fig.5a, tAPE has the highest rank in terms of accuracy and is significantly better than other absolute position encodings due to effectively utilising embedding space to provide an isotropic encoding while holding the distance awareness property. As expected, the model without position encoding

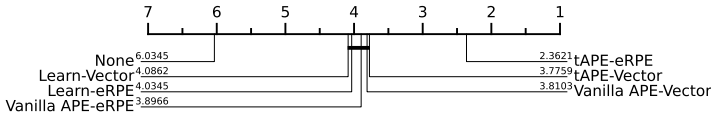


Fig. 6 The average rank of various combination of absolute and relative position encodings.

has the least accurate results, highlighting the importance of absolute position encoding in time series classification. The vanilla APE also improves overall performance despite not being significantly accurate than Learn APE since it has fewer parameters.

Fig.5b shows the critical difference diagram of a single-layer transformer with different relative position encodings. As shown in this figure, eRPE has the highest rank and is significantly better than other encodings in terms of accuracy as it has less number of parameters which is less likely to overfit. It is not surprising that the model without position encoding has the least accurate results, highlighting the importance of relative position encoding and the translation equality property in time series classification. The input-dependent Vector encoding also improves overall performance and is significantly better than None model. Fig.6 shows the critical difference diagram for the various combinations of absolute and relative position encodings. As depicted in this figure, the combination of our proposed tAPE and eRPE is significantly more accurate than all other combinations. This shows the high potential of our encoding methods to incorporate position information into transformers. The combination of Learn and Vector has the least accurate results, most likely due to the high number of parameters.

5.5 Comparing with State-of-the-Art Deep Learning Models

We compare our ConvTran with the following convolution-based and transformer-based models for MTSC:

FCN: Fully Convolutional Neural network is one of the most accurate deep neural networks for MTSC [4] reported in the literature.

ResNet: Residual Network is also one of most accurate deep neural networks for both univariate TSC and MTSC[4] reported in the literature.

Disjoint-CNN: One of the accurate and lightweight CNN-based models that factorize convolution kernels into disjoint temporal and spatial convolutions [26].

Inception-Time: The most accurate deep learning univariate TSC and MTSC algorithm to date. [5, 21].

TST: A transformer-based model for MTSC [12].

Fig. 7 shows the average rank of ConvTran on 32 MTS datasets againsts all convolutional-based and/or transformer-based methods. This figure shows that on average, ConvTran has the lowest average rank and is more accurate than all

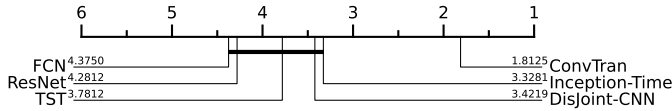


Fig. 7 The average rank of ConvTran against all deep learning based methods on all 32 MTS datasets.

other methods. It is important to observe that ConvTran is significantly more accurate than its predecessors, i.e., a convolution based model, Disjoint-CNN as well as the transformer based model, TST. This indicates the effectiveness of adding tAPE and eRPE to transformers. Table 2 presents the classification accuracy of each method on all 32 datasets and the highest accuracy for each dataset is highlighted in bold. In this table datasets are sorted based on the number of training samples per class. Considering Fig. 7 and Table 2 we can conclude that ConvTran is the most accurate TSC method on average on all 32 benchmark datasets and particularly has superior performance in datasets in which there are enough data to train (i.e., the number of training samples per class is more than 100) and wins on all 12 datasets except one.

5.6 Benchmark against State-of-the-Art Models

Given the experiments on the 32 datasets show that our ConvTran model has the best performance compared to all the other convolution and transformers based models, we now proceed to benchmark it against the state-of-the-art MTSC models, i.e., both deep learning and non-deep learning models. We compare HC2, CIF and ROCKET models on only 26 out of 32 MTSC benchmarking datasets [5] because the other six datasets are either large in terms of training sample or have varied series lengths that make it almost impossible to run HC2 on them. For having detailed insights into the ConvTran performance we provide a pair-wise comparison between our proposed model and each of these models.

As shown in Fig. 8 our proposed model mostly outperforms HC2, ROCKET, CIF, and Inception-Time on the datasets with 100 or more training samples per class (marked with a blue circle). However, state-of-the-art models outperform ConvTran on datasets with few training instances such as **EigenWorms** with 26 train sample per-class. Indeed, as shown in Table 2, all CNN based models fail to perform competitively on the **EigenWorms** dataset. Note that ConvTran is the most accurate among all CNNs on this dataset. This is due to the limitation of CNN-based models, which cannot capture long-term dependencies in the high length time series. Adding a transformer improves the performance, but it still requires more training samples to perform as well as other models.

It is also interesting to observe from Figs. 8a and 8c that HC2 and CIF perform better than ConvTran on the **EthanolConcentration** dataset. Considering that this dataset is based on spectra of water-and-ethanol, hence

Table 2 Average accuracy of six deep learning based models over 32 multivariate time series datasets. Datasets are sorted based on the number of training samples per-class. The highest accuracy for each dataset is highlighted in bold.

DataSets	Avg Train	ConvTran	TST	IT	Disjoint-CNN	FCN	ResNet
Ford	17300	0.7805	0.7655	0.7628	0.7422	0.6353	0.687
HAR	8400	0.9098	0.8831	0.8775	0.8807	0.8445	0.8711
FaceDetection	2945	0.6722	0.6542	0.5885	0.5665	0.5037	0.5948
Insectwingbeat	2500	0.7132	0.6748	0.6956	0.6308	0.6004	0.65
PenDigits	750	0.9871	0.9694	0.9797	0.9708	0.9857	0.9771
ArabicDigits	660	0.9945	0.9749	0.9872	0.9859	0.9836	0.9832
LSST	176	0.6156	0.2846	0.4456	0.5559	0.5616	0.5725
FingerMovement	158	0.56	0.58	0.56	0.54	0.53	0.54
MotorImagery	139	0.56	0.48	0.53	0.49	0.55	0.52
SelfRegSCP1	134	0.918	0.86	0.8634	0.8839	0.7816	0.8362
Heartbeat	102	0.7853	0.6975	0.6248	0.717	0.678	0.7268
SelfRegSCP2	100	0.5833	0.5333	0.4722	0.5166	0.4667	0.5
PhonemeSpectra	85	0.3062	0.089	0.1586	0.2821	0.1599	0.1596
CharacterTraject	72	0.9922	0.9825	0.9881	0.9945	0.9868	0.9945
EthanolConcen	66	0.3612	0.151	0.3489	0.2775	0.3232	0.3155
HandMovement	40	0.4054	0.5405	0.3783	0.5405	0.2973	0.2838
PEMS-SF	39	0.8284	0.7572	0.8901	0.8901	0.8324	0.7399
RacketSports	38	0.8618	0.8815	0.8223	0.8355	0.8223	0.8223
Epilepsy	35	0.9855	0.9492	0.9928	0.8898	0.9928	0.9928
JapaneseVowels	30	0.9891	0.9837	0.9702	0.9756	0.973	0.9135
NATOPS	30	0.9444	0.95	0.9166	0.9277	0.8778	0.8944
EigenWorms	26	0.5934	0.4503	0.5267	0.5934	0.4198	0.4198
UWaveGesture	15	0.8906	0.8906	0.9093	0.8906	0.85	0.85
Libras	12	0.9277	0.8222	0.8722	0.8577	0.85	0.8389
ArticulatoryWord	11	0.9833	0.9833	0.9866	0.9866	0.98	0.98
BasicMotions	10	1	0.975	1	1	1	1
DuckDuckGeese	10	0.62	0.5	0.36	0.5	0.36	0.24
Cricket	9	1	1	0.9861	0.9772	0.9306	0.9722
Handwriting	6	0.3752	0.2752	0.3011	0.2372	0.376	0.18
ERing	6	0.9629	0.9296	0.9296	0.9111	0.9037	0.9296
AtrialFibrillation	5	0.4	0.2	0.2	0.4	0.3333	0.3333
StandWalkJump	4	0.3333	0.3333	0.4	0.3333	0.4	0.4

interval and shapelet-based approaches which are also components of HC2 perform better. On the other hand, ROCKET has a few wins compared to ConvTran (Fig 8b). Most of these datasets where ROCKET performs better, such as the `StandWalkjump` dataset have a small number of time series instances per class. For instance, `StandWalkjump` has 3 classes with 12 training instances, which is 4 time series per class. This is insufficient to train large number of parameters in deep learning models such as ConvTran to achieve better performance. Note, as mentioned, these results are for 26 datasets only, excluding six datasets for which we could not run HC2 (which has high computational complexity and is limited to be applied on variable-length time series). Among excluded datasets, 4 of them are large datasets from which ConvTran could have benefited. Considering this, ConvTran still achieves competitive performance compared to SOTA deep and non-deep models.

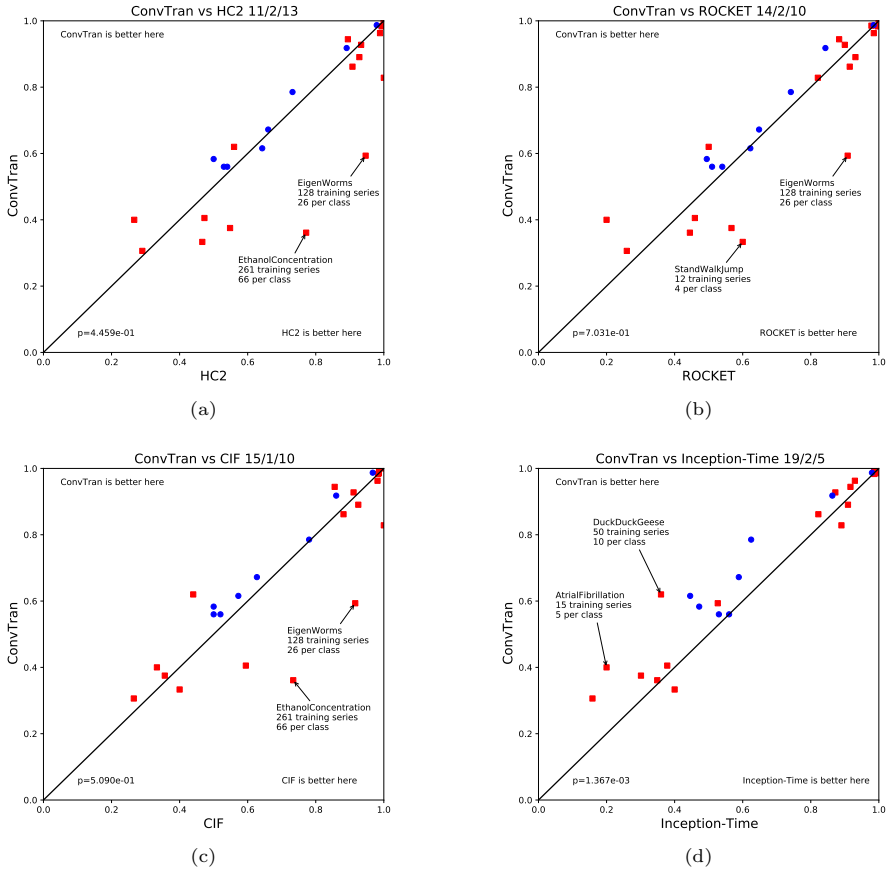


Fig. 8 Pairwise comparison of ConvTran with the state of the art models: (a) HC2, (b) ROCKET, (c) CIF (d) and Inception-Time. The datasets with 100 training samples per class or more are marked with a blue circle, while the others are marked with a red square. The three values at the top of each figure show the number of win/draw/loss from left to right

5.7 ConvTran vs ROCKET Efficiency and Effectiveness

To provide further insight into the efficiency of our model on datasets of varying sizes, we conducted additional experiments on the largest UEA dataset *InsectWingBeat* with 25,000 series for training. We compare the training time and test accuracy of our proposed ConvTran and ROCKET on random subsets of 5,000, 10,000, 15,000, 20,000, and 25,000 training samples.

The results depicted in Figure 9 demonstrate that ROCKET has faster training time than ConvTran on smaller datasets, specifically on the 5k and 10k datasets while achieving similar training time to ConvTran on the 15k set. However, our deep learning-based model, ConvTran, demonstrates faster training times with increasing data quantity, as expected. Additionally, we also observed from the figure that ConvTran is consistently more accurate than

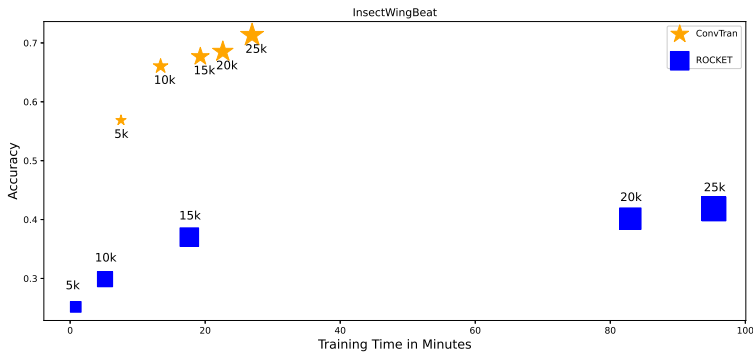


Fig. 9 Comparison of runtime and accuracy between ConvTran and ROCKET on UEA largest dataset InsectWingBeat with 25,000 training samples. The figure shows the runtime of the two models on datasets with different sizes, and their corresponding classification accuracy.

ROCKET on this dataset. We refer interested readers to Appendix A.1 for a more comprehensive exploration of the empirical evaluation of efficiency and effectiveness on all datasets. Notably, ConvTran demonstrates faster inference time compared to ROCKET across all datasets. It is important to note that all the ConvTran experiments are performed on GPUs, whereas ROCKET experiments are performed on a CPU (please refer to Section 5 for computing system details).

6 Conclusion

This paper studies the importance of position encoding for time series for the first time and reviews existing absolute and relative position encoding methods in time series classification. Based on the limitations of the current position encodings for time series, we proposed two novel absolute and relative position encodings specifically for time series called tAPE and eRPE, respectively. We then integrated our two proposed position encodings into a transformer block and combine them with a convolution layer and presented a novel deep-learning framework for multivariate time series classification (ConvTran). Extensive experiments show that ConvTran benefits from the position information, achieving state-of-the-art performance on Multivariate time series classification in deep learning literature. In future, we will study the effectiveness of our new transformer block in other transformer-based TSC models and other down stream tasks such as anomaly detection.

7 Declarations

Conflict of interest statement: The authors have no competing interests to declare that are relevant to the content of this article.

References

- [1] Lockhart, J.W., Weiss, G.M., Xue, J.C., Gallagher, S.T., Grosner, A.B., Pulickal, T.T.: Design considerations for the wisdm smart phone-based sensor mining architecture. In: *International Workshop on Knowledge Discovery from Sensor Data*, pp. 25–33 (2011)
- [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] Bagnall, A., Lines, J., Bostrom, A., Large, J., Keogh, E.: The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery* **31**(3), 606–660 (2017)
- [4] 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)
- [5] Ruiz, A.P., Flynn, M., Large, J., Middlehurst, M., Bagnall, A.: The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, 1–49 (2020)
- [6] Dai, Z., Liu, H., Le, Q.V., Tan, M.: Coatnet: Marrying convolution and attention for all data sizes. *Advances in Neural Information Processing Systems* **34**, 3965–3977 (2021)
- [7] Karim, F., Majumdar, S., Darabi, H., Harford, S.: Multivariate lstm-fcns for time series classification. *Neural Networks* **116**, 237–245 (2019)
- [8] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in Neural Information Processing Systems* **30** (2017)
- [9] Hao, Y., Cao, H.: A new attention mechanism to classify multivariate time series. In: *International Joint Conference on Artificial Intelligence* (2020)
- [10] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
- [11] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)

- [12] Zerveas, G., Jayaraman, S., Patel, D., Bhamidipaty, A., Eickhoff, C.: A transformer-based framework for multivariate time series representation learning. In: SIGKDD Conference on Knowledge Discovery & Data Mining, pp. 2114–2124 (2021)
- [13] Kostas, D., Aroca-Ouellette, S., Rudzicz, F.: Bendr: using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers in Human Neuroscience* **15** (2021)
- [14] Shaw, P., Uszkoreit, J., Vaswani, A.: Self-attention with relative position representations. arXiv preprint arXiv:1803.02155 (2018)
- [15] Huang, C.-Z.A., Vaswani, A., Uszkoreit, J., Shazeer, N., Simon, I., Hawthorne, C., Dai, A.M., Hoffman, M.D., Dinculescu, M., Eck, D.: Music transformer. arXiv preprint arXiv:1809.04281 (2018)
- [16] Dufter, P., Schmitt, M., Schütze, H.: Position information in transformers: An overview. *Computational Linguistics* **48**(3), 733–763 (2022)
- [17] Foumani, N.M., Miller, L., Tan, C.W., Webb, G.I., Forestier, G., Salehi, M.: Deep learning for time series classification and extrinsic regression: A current survey. arXiv preprint arXiv:2302.02515 (2023)
- [18] Dempster, A., Petitjean, F., Webb, G.I.: Rocket: exceptionally fast and accurate time series classification using random convolutional kernels. *Data Mining and Knowledge Discovery* **34**(5), 1454–1495 (2020)
- [19] Bagnall, A., Flynn, M., Large, J., Lines, J., Middlehurst, M.: On the usage and performance of the hierarchical vote collective of transformation-based ensembles version 1.0 (hive-cote v1. 0). In: International Workshop on Advanced Analytics and Learning on Temporal Data, pp. 3–18 (2020)
- [20] Middlehurst, M., Large, J., Bagnall, A.: The canonical interval forest (cif) classifier for time series classification. In: 2020 IEEE International Conference on Big Data, pp. 188–195 (2020)
- [21] Fawaz, H.I., Lucas, B., Forestier, G., Pelletier, C., Schmidt, D.F., Weber, J., Webb, G.I., Idoumghar, L., Muller, P.-A., Petitjean, F.: Inceptiontime: Finding alexnet for time series classification. *Data Mining and Knowledge Discovery* **34**(6), 1936–1962 (2020)
- [22] Dempster, A., Schmidt, D.F., Webb, G.I.: Minirocket: A very fast (almost) deterministic transform for time series classification. In: SIGKDD Conference on Knowledge Discovery & Data Mining, pp. 248–257 (2021)
- [23] Tan, C.W., Dempster, A., Bergmeir, C., Webb, G.I.: Multirocket: Effective summary statistics for convolutional outputs in time series classification.

arXiv e-prints, 2102 (2021)

- [24] Middlehurst, M., Large, J., Flynn, M., Lines, J., Bostrom, A., Bagnall, A.: Hive-cote 2.0: a new meta ensemble for time series classification. *Machine Learning* **110**(11), 3211–3243 (2021)
- [25] 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, pp. 1578–1585 (2017)
- [26] Foumani, S.N.M., Tan, C.W., Salehi, M.: Disjoint-cnn for multivariate time series classification. In: 2021 International Conference on Data Mining Workshops, pp. 760–769 (2021)
- [27] Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.-C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 4510–4520 (2018)
- [28] 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)
- [29] Yang, C.-H.H., Tsai, Y.-Y., Chen, P.-Y.: Voice2series: Reprogramming acoustic models for time series classification. In: International Conference on Machine Learning, pp. 11808–11819 (2021)
- [30] Luong, M.-T., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. arXiv preprint arXiv:1508.04025 (2015)
- [31] Huang, Z., Liang, D., Xu, P., Xiang, B.: Improve transformer models with better relative position embeddings. arXiv preprint arXiv:2009.13658 (2020)
- [32] Wu, K., Peng, H., Chen, M., Fu, J., Chao, H.: Rethinking and improving relative position encoding for vision transformer. In: IEEE/CVF International Conference on Computer Vision, pp. 10033–10041 (2021)
- [33] Liang, Y., Cao, R., Zheng, J., Ren, J., Gao, L.: Learning to remove: Towards isotropic pre-trained bert embedding. In: International Conference on Artificial Neural Networks, pp. 448–459 (2021)
- [34] Demšar, J.: Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research* **7**, 1–30 (2006)

Table A1 Comparison of runtime and accuracy between ConvTran and ROCKET on 32 datasets of varying sizes. To facilitate easy identification, superior performance in both accuracy and runtime is highlighted in bold in the table. For a detailed comparison, the runtimes are shown in seconds.

Datasets	Train size	ROCKET			ConvTran		
		Accuracy	Train time	Test time	Accuracy	Train time	Test time
HAR	41546	0.8293	5366.34	11.51	0.9098	2367.77	1.82
Ford	28839	0.6051	6863.81	11.91	0.7805	1619.42	0.95
InsectWingbeat	25000	0.4182	5721.04	41.5	0.7132	1617.82	5.47
PenDigits	7494	0.984	65.26	0.99	0.9871	401.1	0.59
ArabicDigits	6599	0.9932	75.59	10.38	0.9945	376.7	0.37
FaceDetection	5890	0.5624	53.23	11.99	0.6722	413.39	0.83
PhonemeSpectra	3315	0.1894	42	37.22	0.3062	202.27	0.89
LSST	2459	0.5251	5.84	3.52	0.6156	148.07	0.48
CharacterTrajec	1422	0.9916	8.4	7.72	0.9922	89.61	0.28
FingerMovement	316	0.55	0.96	0.35	0.56	21.33	0.02
MotorImagery	278	0.56	45.39	16.26	0.56	386	0.81
ArticularyWord	275	0.9933	2.09	2.19	0.9833	19.76	0.08
JapaneseVowels	270	0.9568	0.57	0.67	0.9891	20.6	0.13
SelfRegSCP1	268	0.8601	10.2	11.25	0.918	45.54	0.27
PEMS-SF	267	0.8266	3.53	2.13	0.8284	28.08	0.09
EthanolConcen	261	0.4448	14.59	14.32	0.3612	131.58	0.69
Heartbeat	204	0.7414	4.57	4.59	0.7853	17.13	0.09
SelfRegSCP2	200	0.5833	10.78	9.65	0.5833	50.05	0.22
NATOPS	180	0.8944	0.6	0.58	0.9444	14.61	0.04
Libras	180	0.8667	0.36	0.29	0.9277	11.51	0.04
HandMovement	160	0.4189	3.31	1.7	0.4054	11.29	0.03
RacketSports	151	0.9078	0.29	0.32	0.8618	11.86	0.03
Handwriting	150	0.5376	0.81	3.92	0.3752	11.85	0.23
Epilepsy	137	0.971	0.91	0.93	0.9855	10.52	0.03
EigenWorms	128	0.8702	107.48	111.42	0.5934	225.71	0.7
UWaveGesture	120	0.9188	1.21	3.07	0.8906	10.2	0.09
Cricket	108	1	5.79	4.07	1	32.1	0.1
DuckDuckGeese	50	0.5	1.59	1.76	0.62	9.46	0.05
BasicMotions	40	1	0.25	0.27	1	4.45	0.01
ERing	30	0.9851	0.17	0.7	0.9629	3.17	0.06
AtrialFibrillation	15	0.2	0.39	0.41	0.4	1.99	0.01
StandWalkJump	12	0.5333	1.52	1.65	0.3333	14.56	0.09

Appendix A

A.1 Empirical Evaluation of Efficiency and Effectiveness

The results presented in Table A1 demonstrate that ConvTran outperforms ROCKET in terms of both train time and test accuracy on larger datasets with more than 10k samples. However, ROCKET has a better train time on smaller datasets. Nevertheless, even on small datasets, ConvTran achieves acceptable accuracy within a reasonable train time. It is worth noting that the performance of ConvTran improves as the dataset size increases, indicating that our model is suitable for scaling to larger datasets.

A.2 ConvTran vs non-deep learning SOTA Models

Table A2 compares the performance of ConvTran against three non-deep learning models - ROCKET, HC2, and CIF - on different datasets with varying training sample sizes. The table presents the accuracy of each model on each dataset, with boldface indicating superior accuracy. “-” denotes non-runnable methods, either due to computation complexity or inability to handle various length series.

Overall, ConvTran outperforms the non-deep learning models on 19 out of 32 datasets (for the HC2 and CIF models, we only have results for 26 datasets, and ConvTran outperforms the other models in 13 out of the 26). It performs better on datasets with larger training sample sizes, such as InsectWingBeat, while other models perform better on datasets with fewer training samples, such as StandWalkJump, which only has 12 training samples. Additionally, the table shows that some of the non-deep learning models failed to handle specific datasets due to either computational complexity or the inability to handle varying input series lengths. For example, we were not able to run HC2 and CIF on the larger HAR, Ford, and InsectWingbeat datasets due to computational complexity. They were also not designed to handle varying length time series such as the CharacterTrajectories, SpokenArabicDigits, and JapaneseVowels datasets.

Table A2 Comparison of ConvTran and Non-Deep Learning Models (ROCKET, HC2, CIF) on Varying Training Sample Sizes. Bold Face Font Indicates Superior Accuracy, ‘-’ Denotes Non-Runnable Methods due to Computation Complexity or Inability to Handle Various Length Series.

Datasets	Train Size	ConvTran	ROCKET	HC2	CIF
HAR	41546	0.9098	0.8293	-	-
Ford	28839	0.7805	0.6051	-	-
InsectWingbeat	25000	0.7132	0.4182	-	-
PenDigits	7494	0.9871	0.984	0.9791	0.9674
SpokenArabicDigits	6599	0.9945	0.9932	-	-
FaceDetection	5890	0.6722	0.5624	0.6603	0.6271
PhonemeSpectra	3315	0.3062	0.1894	0.2905	0.2654
LSST	2459	0.6156	0.5251	0.6427	0.5726
CharacterTrajectories	1422	0.9922	0.9916	-	-
FingerMovements	316	0.56	0.55	0.53	0.52
MotorImagery	278	0.56	0.56	0.54	0.5
ArticulatoryWordRecognition	275	0.9833	0.9933	0.9933	0.9833
JapaneseVowels	270	0.9891	0.9568	-	-
SelfRegulationSCP1	268	0.918	0.8601	0.8908	0.8601
PEMS-SF	267	0.8284	0.8266	1	1
EthanolConcentration	261	0.3612	0.346	0.7719	0.7338
Heartbeat	204	0.7853	0.678	0.7317	0.7805
SelfRegulationSCP2	200	0.5833	0.5833	0.5	0.5
NATOPS	180	0.9444	0.8944	0.8944	0.8556
Libras	180	0.9277	0.8667	0.9333	0.9111
HandMovementDirection	160	0.4054	0.4189	0.473	0.5946
RacketSports	151	0.8618	0.9078	0.9078	0.8816
Handwriting	150	0.3752	0.5376	0.5482	0.3565
Epilepsy	137	0.9855	0.971	1	0.9855
EigenWorms	128	0.5934	0.8702	0.9466	0.916
UWaveGestureLibrary	120	0.8906	0.9188	0.9281	0.925
Cricket	108	1	1	1	0.9861
DuckDuckGeese	50	0.62	0.5	0.56	0.44
BasicMotions	40	1	1	1	1
ERing	30	0.9629	0.9593	0.9889	0.9815
AtrialFibrillation	15	0.4	0.1333	0.2667	0.3333
StandWalkJump	12	0.3333	0.5333	0.4667	0.4
Wins or Draw	-	19	7	13	4