# rTsfNet: a DNN model with Multi-head 3D Rotation and Time Series Feature Extraction for IMU-based Human Activity Recognition

YU ENOKIBORI, Graduate School of Informatics, Nagoya University, Japan

Although many deep learning (DL) algorithms have been proposed for the IMU-based HAR domain, traditional machine learning that utilizes handcrafted time series features (TSFs) still often performs well. It is not rare that combinations among DL and TSFs show better accuracy than DL-only approaches. However, there is a problem with time series features in IMU-based HAR. The amount of derived features can vary greatly depending on the method used to select the 3D basis. Fortunately, DL's strengths include capturing the features of input data and adaptively deriving parameters. Thus, as a new DNN model for IMU-based human activity recognition (HAR), this paper proposes rTsfNet, a DNN model with Multi-head 3D Rotation and Time Series Feature Extraction. rTsfNet automatically selects 3D bases from which features should be derived by extracting 3D rotation parameters within the DNN. Then, time series features (TSFs), based on many researchers' wisdom, are derived to achieve HAR using MLP. Although rTsfNet is a model that does not use CNN, it achieved higher accuracy than existing models under well-managed benchmark conditions and multiple datasets: UCI HAR, PAMAP2, Daphnet, and OPPORTUNITY, all of which target different activities.

CCS Concepts: • **Computing methodologies → Machine learning algorithms**; • **Human-centered computing → Ubiquitous and mobile computing**.

Additional Key Words and Phrases: Human activity recognition, Time series feature, Deep neural network, Algorithm, Benchmark setup, Multi-head 3D rotation, IMU

## 1 INTRODUCTION

Since the advent of deep learning (DL), many DL methods have been proposed for IMU-based human activity recognition (HAR). Many related mechanisms also have been proposed frequently, such as HAR dataset federation using DNN [14], self-supervised learning for HAR [9, 32], and so on. Focusing on their core recognition part, many of them are influenced by the image processing field and apply CNN to IMU output values or RNN to CNN values to derive features and discriminative results without handcrafted time series features.

However, in the IMU-based HAR domain, traditional machine learning (ML) that utilizes handcrafted time series features, which are comprised of the intellectual wisdom of many researchers, still often performs well. For example, as shown in Table 1, the reference result of the UCI HAR [2] Benchmark setup that uses SVM and 561 handcrafted features has an accuracy of 96.37, a value that exceeds almost DL-based approaches. The second-ranked ML is that of C.A. Ronao et al. [28] whose approach uses Hierarchical Continuous HMMs. Their approach shows an accuracy of

Author's address: Yu Enokibori, enokibori@i.nagoya-u.ac.jp, Graduate School of Informatics, Nagoya University, #361, IB-south-tower, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi, Japan, 464-8603.

Table 1. Comparition on UCI HAR

| study | year | note | acc | mf1 | wf1 | note |
|---|---|---|---|---|---|---|
| Y. Enokibori | - | proposed | 97.76 | 0.9779 | 0.9776 | - |
| A. Ignatov [12] | 2018 | CNN + stat. f. + data-centering | 97.63 | - | - | use test set info. |
| W. Jiang et al.[13] | 2015 | 2D FFT + CNN + SVM | 97.59 | - | - | ensemble |
| Y. Enokibori | - | proposed, w/o mh-3D-rotation | 96.71 | 0.9675 | 0.9670 | - |
| L. Lu et al. [20] | 2022 | CNN + GRU | 96.67 | - | - | - |
| D. Anguita et al.[2] | 2013 | SVM + 561 hc. features | 96.37 | 0.9638 | 0.9636 | reference |
| A. Ignatov [12] | 2018 | CNN + stat. feature | 96.06 | - | - | - |
| C.A. Ronao et al.[27] | 2016 | FFT + CNN | 95.75 | - | - | - |
| Z.N. Khan et al.[15] | 2021 | Attention + CNN | 95.40 | - | - | - |
| A. Dahou et al.[6] | 2022 | DNN + BAOA + SVM | 95.33 | 0.9533 | - | ensemble |
| Y. Dong et al.[8] | 2020 | DSMT | 95.31 | 0.9467 | - | - |
| W. Jiang et al.[13] | 2015 | 2D FFT + CNN | 95.18 | - | - | - |
| K. Xia et al.[37] | 2020 | LSTM-CNN | 95.14 | 0.9525 | - | - |
| J. He et al.[10] | 2019 | SVL + RNN + Attention | 94.80 | - | - | - |
| C.A. Ronao et al.[27] | 2016 | CNN | 94.79 | - | - | - |
| C. Xu et al.[38] | 2019 | CNN + LSTM | 94.60 | - | - | - |
| N. Tufre et al.[34] | 2020 | LSTM | 93.70 | 0.8260 | - | - |
| Y. Zhao et al.[39] | 2018 | Bi-LSTM | 93.60 | - | - | - |
| K. Wang et al.[36] | 2019 | Attention + CNN | 93.41 | - | - | - |
| C.A Ronao[28] | 2017 | Hierarchical Continuous HMM | 93.18 | - | - | - |
| V. Bianchi et al.[5] | 2019 | CNN | 92.50 | 0.9210 | | - |
| Y. Li et al.[19] | 2014 | Stacked Auto Encoders + SVM | 92.16 | - | - | - |
| Y. Li et al.[19] | 2014 | PCA + SVM | 91.82 | - | - | - |
| C.A. Ronao et al.[25] | 2014 | HMM | 91.76 | - | - | - |
| H. Li et al.[18] | 2020 | Bi-LSTM | 91.21 | 0.9001 | - | - |
| C.A. Ronao et al.[26] | 2015 | CNN | 90.89 | - | - | - |
| S. Seto et al.[29] | 2015 | DTW | 89.00 | - | - | - |
| D. Anguita et al.[1] | 2012 | Handcrafted features + SVM | 89.00 | - | - | - |
| Y.J. KIM et al.[16] | 2015 | HMM | 83.51 | - | - | - |

93.18 under the UCI HAR Benchmark setup, which is only 2 points lower than the DL of 95.40 with CNN and Attention by Z.N. Khan et al. [15] In addition, methods that combine CNN and handcraft features, such as FFT, as proposed by A. Ignatiov [12], W. Jiang [13], and so on., achieved higher accuracy than methods using the Attention mechanism, which is one of the newest DL approaches. The above facts indicate that handcraft-based time series features (TSF) may suit IMU-based HAR.

On the other hand, there is a problem with time series features in IMU-based HAR. The amount of derived features can vary greatly depending on the method used to select the 3D basis. For example, even if a sensor is mounted at the same wrist position, different features will be derived if the mounting angle is different. This problem is commonly solved by extracting an effective 3D basis using PCA. The use of the L2 norm is also often used to ignore those basis differences. However, selecting only one 3D basis by PCA is undeniably inadequate for HAR, which is conducted under various conditions. In addition, this approach is not good at selecting axes by comprehensively considering multiple

sensors, such as accelerometers, gyroscopes, and magnetometers. Only using the L2 norm tends to reduce the amount of information, which in turn lowers the final accuracy limit.

Fortunately, DL's strengths include capturing the features of input data and adaptively deriving parameters. Thus, this paper proposes rTsfNet, which selects multiple 3D bases in a DNN and derives time series features from the data in them. Although rTsfNet is a model that does not use CNNs, its accuracy still outperforms existing models under multiple benchmark setups.

The remainder of this paper is organized as follows: Section 2 describes the verifiability issue of IMU-based HAR and explores related studies based on well-managed benchmark setups. Our proposed model, rTsfNet, is discussed in Section 3. Section 4 evaluates and compares its performance with other studies. After rTsfNet's potential is described in Section 5, this paper's conclusion is presented in Section 6.

## 2 DISCUSSION ABOUT IMU-BASED HAR BENCHMARK AND STUDIES

This section summarizes the current status of IMU-based HAR. First, we discuss the verifiability issue as a prerequisite for comparing IMU-based HAR studies. Then, we summarize related studies based on the UCI HAR Benchmark setup, which has a clean evaluation setup, allows direct comparison, and has been used in many studies. Note that this paper does not discuss ensemble learning, data augmentation, semi-supervised learning, and other techniques that extend basic algorithms.

### 2.1 Verifiability issue

*2.1.1 Direct comparability issue with benchmark setup.* Many datasets for IMU-based human action recognition have been proposed [2, 3, 17, 22, 23, 31]. Many IMU-based HAR studies have evaluated their methods and compared them with other methods on identical datasets. However, more than just using the same dataset is required to properly compare action recognition methods properly. The segmentation and train/test set split methods must be unified.

For the segmentation method, the window and sliding sizes must be clear, at least for the test set. For example, the difficulty of detecting a gait activity with 0.5 seconds of window size data differs from detecting such activity with 10 seconds.

The train/test set split must be subject-based, trial-based, or time-based. Simple random ratio-based splitting is not acceptable in IMU-based HAR because the data at times $t$ and $t + 1$ are similar. If they belong to separate train/test sets, overfitting will occur. The accuracy in the simple random ratio-based splitting test set is almost the same as in the train set.

Moreover, such notations as "a ratio-based split based on subject" are insufficient because the accuracy varies greatly depending on how the test set is selected. For example, in the UCI HAR dataset, subject 14 is an anomaly, and the accuracy changes depending on whether subject 14 is included in the test set. If the train/test set split is unclear, an overfit result can be generated by including the anomaly in the train set. Train and test sets, and also a validation set if possible, must be defined clearly, such as by subject IDs, trial IDs, and file names.

Few studies in the field of IMU-based HAR field satisfy these perspectives and allow direct comparisons. One rare exception where multiple studies can be directly compared is a group of studies using the benchmark setup of the UCI HAR dataset. This dataset is provided in a pre-segmented form, and the train/test set split is subject-based and clearly defined in the subject ID.

Another valid comparison environment is the iSPLInception Benchmark setup. iSPLInception Benchmark setup defines the segmentation settings, and train/test set splits as subject-based, trial-based, or file-based for UCI HAR,

PAMAP2[22], OPPORTUNITY[23], and Daphnet[3], considering class sample ratios. Most importantly, the source code of the data handling is opened.

Therefore, as direct comparisons, we use studies that are basically evaluated in the UCI HAR. Studies evaluated in the iSPLInception Benchmark setup will also be compared. However, as described later in Subsection 2.3, the iSPLInception Benchmark setup has problems with time warping and dirty segmentation, both of which have a small impact on PAMAP2 and Daphnet but a massive impact on OPPORTUNITY. At least for OPPORTUNITY, the iSPLInception Benchmark setup should not be used. Therefore, in this paper, we define a new benchmark setup, called the IMU-based HAR Benchmark, and make it open, to solve these problems and to help future studies.

*2.1.2 LOSO CV.* As for the train/test set split, leave-one-subject-out cross-validation (LOSO CV) is also valid since it clearly determines the train/test set split. However, unfortunately, segmentation is often unclear/different in many studies, and so direct comparison is difficult.

Another problem with LOSO CV is that it increases both the parameter search and evaluation times. Currently, most datasets in the IMU-based HAR domain consist of about ten subjects, an amount that rises up to 50 for a big dataset, in the IMU-based HAR domain. However, this is very small compared to image datasets. If possible, the datasets of the IMU-based HAR domain should be comprised of 1,000 or 10,000 subjects to consider the individual differences. This idea is one of the future works of the IMU-based HAR domain. However, when the dataset size increases, the LOSO CV will cause an evaluation time problem. If increases of evaluation times are acceptable, it is better to evaluate the generality of methods using different datasets that have various types of tasks. In other words, such evaluation should be done on multiple datasets in a well-designed benchmark setup. Therefore, in this paper, we did not select the LOSO CV for our evaluation setup.

*2.1.3 Re-verifiability.* Unfortunately, unlike the image recognition domain, many IMU-based HARs are closed sources and cannot be re-verified. Only a few can be verified, such as DeepConvLSTM and iSPLInception. This situation is unhealthy. For healthier evaluations, the results should be eliminated that others cannot replicate. Such re-verifiability is especially important when identification accuracy approaches 100% and when tiny improvements are being compared.

*2.1.4 Metrics.* Many studies use accuracy, F1 scores, and weighted F1 scores (wf1) as representative values for the performance of other classifiers. However, F1 scores do not match well for multiclass recognition. Wf1's value is almost the same value as its accuracy; it is not a very effective value. Neither should accuracy nor wf1 should not be used, especially for extremely class-imbalanced datasets. For example, the iSPLInception Benchmark setup for the Daphnet dataset has 2103 samples for class 1 and 126 for class 2. Simply answering every sample as class 1 yields an accuracy and wf1 of 92.84% and 0.9284. Such high values are misleading. On the other hand, the macro F1 score (mf1) shows 0.5 in the same situation. If a representative value were chosen, mf1 is preferable. If the imbalance is not too extreme, then the accuracy comparison is acceptable. However, in most cases, mf1 shows similar values of accuracy.[1]

In the datasets discussed so far, OPPORTUNITY, which excludes the null labels, UCI HAR, and PAMAP2 are class-imbalanced but not excessively. Therefore, accuracy, which is used in many studies, can be used for direct comparison. On the other hand, mf1 should be used for Daphnet.

---

[1]They can be recalculated if the confusion matrix is provided in terms of numbers rather than percentages. Confusion matrices presented as percentages are only valid on one axis. In addition, since this approach basically hides information, it is detrimental and must not be used unless the class is balanced.

*2.1.5 Summary.* Based on the above discussion, this study mainly compares and discusses studies using the UCI HAR benchmark setup. We also compare the results of several algorithms evaluated with iSPLInception-like benchmark setups called IMU-based HAR Benchmark setups.

## 2.2 IMU-based HAR evaluated by UCI HAR

As shown above, Table 1 summarizes the accuracy of the studies that use the UCI HAR benchmark setup, including the rTsfNet results.

Until around 2015, IMU-based HARs were dominated by handcrafted features, such as statistical features, and traditional ML, such as HMM and SVM. For example, C.A Ronao et al. [28] achieved 93.18% accuracy in the UCI HAR benchmark setup with a Hierarchical Continuous HMM. D. Anguita et al. [2] showed 96.37% accuracy as the reference performance of the UCI HAR dataset with SVM and 561 handcrafted features.

With the subsequent significant development of DNNs in the image recognition domain, IMU-based HAR has also become DNN-based. C.A. Ronao et al. [27] achieved 94.79% accuracy in the UCI HAR benchmark setup with CNN. Some studies have been combining CNN and RNN. K. Xia et al. [37] achieved 95.14% accuracy with a combination of LSTM and CNN. L. Lu et al. [20] achieved 96.67% accuracy by combining LSTM and GRU with a multichannel stream approach. Some other studies use Attention, which is a key mechanism of the Transformer's great success with LLMs. Z.N. Khan et al.[15] achieved 95.40% accuracy with a combination of CNN and Attention.

Some studies have used DNNs and handcrafted features together. For example, C.A. Ronao et al. [27] achieved 95.75% accuracy by inputting frequency features derived by FFT into a CNN. A. Ignatov [12] achieved 96.06% accuracy by concatenating handcrafted features after feature extraction in a CNN. They also got 97.63% accuracy by performing data centering, but this method is not treated in this paper because it uses information from its test dataset.

The combination of DL, ML, and handcrafted features shows the highest performance among related studies in Table 1. W. Jiang et al. [13] achieved 97.59% accuracy with a combination of FFT, CNN, and SVM.

As described above, the importance of the methods using handcrafted features is clear, since their accuracy surpasses the methods using state-of-the-art DNN mechanisms.

## 2.3 Re-verifiable Works

Unfortunately, the related studies described in the previous subsection are unverifiable, unlike the image recognition domain. DeepConvLSTM and iSPLInception are rare studies whose sources are publicly available and can be re-verified.

DeepConvLSTM is a combined CNN and LSTM model that has been evaluated on the OPPORTUNITY dataset. Benchmark setups, which are also available in the public source code, are clearly defined for both segmentation and train/test set splits. However, using DeepConvLSTM for direct comparison poses several issues that must be confronted.

The first concerns the handling of NaN values. In the benchmark setup of DeepConvLSTM, this NaN value is complemented linearly without any time limit. Therefore, a variation rule can be specifically created in the sensor values, an advantageous result for interpretation by RNNs and other methods.[2]

The second issue is the identification target difference issue. DeepConvLSTM's identification target is the last label of the segment, a strategy that is different from such common tasks as estimating modes. Therefore, we do not use this setup or perform a direct comparison with DeepConvLSTM.[3]

---

[2]The OPPORTUNITY dataset has a continuous and long-lasting high occurrence of NaN values specific to certain behaviors and sensors in the 12 accelerometers.

[3]In addition, when concatenating multiple trials, the segmentation does not consider the breaks, and time warps occur within some segmentations. However, this issue is a tiny amount, and so the impact is minimal.

Table 2. Activity length summary of OPPORTUNITY

| length | % less than N |
|--------|---------------|
| 30 | 1.72 |
| 32 | 2.75 |
| 45 | 10.65 |
| 60 | 26.46 |
| 75 | 54.98 |
| 90 | 76.63 |

The iSPLInception, which is a model that applies the inception [30] structure proposed in image recognition to HAR, is evaluated using the following four datasets: UCI HAR, a dataset of basic behaviors under the single sensor and simple measurement conditions; PAMAP2, a dataset of basic behaviors under multi-sensor and complex conditions; OPPORTUNITY, a dataset of such daily behaviors as opening a door acquired with multi-sensors, and Daphnet, a dataset that targets the freezing phenomenon of Parkinson's disease, whose capture is challenging to capture with accelerometers compared to basic behavior. The published source code includes the benchmark setups, and both segmentation and train/test set splits and the validation set split, are clearly defined.

However, since the iSPLInception Benchmark setup defines a different split method for the UCI HAR than the original train/test set split, it cannot be used for a direct comparison with other studies validated using UCI HAR. In addition, this benchmark setup has time warps in the segmentation due to the simple exclusion of samples containing NaN values and null labels, and no consideration of the breaks among subjects, trials, and files in the datasets except for the UCI HAR.

Segmentation affected by simply excluding samples containing NaN values and null labels in PAMAP2 and Daphnet is also minimal and has little impact. Segmentation involving inter-user and inter-trial breaks is also negligible and has little impact on all the datasets.

Compared to these datasets, the OPPORTUNITY dataset significantly impacts the issue described above. The OPPORTUNITY dataset almost always has samples labeled null between each activity. In addition, as shown in Table 2, a majority of its activities end in fewer than 90 samples, although the iSPLInception Benchmark setup divides OPPORTUNITY into 90 samples. Therefore, 76.63% of the segmentations have at least one time warp. Such dirty segmentation is detrimental to discriminators that consider order, signal frequency, etc.

To solve the issues described above, this study defined an IMU-based HAR Benchmark setup based on the iSPLInception Benchmark setup with an updated NaN value handling and segmentation method, as shown in Subsection 4.7. Although the IMU-based HAR Benchmark and the iSPLInception Benchmark setups are not strictly comparable, the PAMAP2 and Daphnet data sets are less affected by the modifications, and so we use the values before and after the modifications are used for comparison and verification in this paper. On the other hand, since OPPORTUNITY is affected by a large number of modifications, we do not directly compare it. rTsfNet's performance in the iSPLInception Benchmark setup for OPPORTUNITY is presented for comparison, although its advisory values should not be compared in future studies.

Fig. 1. rTsfNet



Fig. 2. MLP Block



Fig. 3. Tsf Mixer sub-Block

## 3 RTSFNET

rTsfNet selects multiple 3D bases in the DNN and derives and extracts time series features from the data in them. An overview is illustrated in Figure 1. The network structure proposed in this paper does not use Residual, SE, Attention, LSTM, ensemble, and so on. We dare to propose it as a basic structure for other networks like CNN. Figures 1 to 6 show the overall picture.

In rTsfNet, the values from 3-axis sensors that can rotate in 3D are treated separately from other sensor values. The values of 3-axis sensors are subjected to multiple rotations in the Multi-head 3D Rotation Block. Then, we concatenate

TSF Mixer Block ($n^{\text{bk}}_{1 \text{ to } 6}$, $n^{\text{d}}_{1 \text{ to } 6}$)

TSF Mixing sub-Block
($n^{\text{bk}}_1$, $n^{\text{d}}_1$, $n^{\text{bk}}_2$, $n^{\text{d}}_2$, $n^{\text{ch}}$)

binary selection
weight for channel

binary selection
weight for axes

TSF Mixing sub-Block
($n^{\text{bk}}_3$, $n^{\text{d}}_3$, $n^{\text{bk}}_4$, $n^{\text{d}}_4$ $n^{\text{ax}}$)

Mixing TSF

in

MLP Block
($n^{\text{bk}}_5$, $n^{\text{d}}_5$)

TSF for each axis

out

MLP Block
($n^{\text{bk}}_6$, $n^{\text{d}}_6$)

Mixing all

Flatten

Fig. 4. Tsf Mixer Block

Multi-head 3D Rotation Block (block_params, $n^{\text{h}}$, $n^{\text{bk}}_{1 \text{ to } 7}$, $n^{\text{d}}_{1 \text{ to } 7}$)

in

L2

Repeat until $n^{\text{h}}$

Rotation Parameter
Calculation Block

Rotation Parameter
Calculation Block ($n^{\text{bk}}_{1 \text{ to } 7}$, $n^{\text{d}}_{1 \text{ to } 7}$)

Tagging     Splitting and extracting TSF with block_params
for each axis of sensors on each block

3D
rotation
matrix

Apply 3D
rotations for each
tri-axes feature

$n^{\text{h}}$

1

out

Fig. 5. Multi-head 3D Rotation Block

the L2 norm of the values of the 3-axis sensors that can rotate 3D, the sensor values after the rotation, and the values of the sensors that cannot rotate 3D and extract a time series feature from each axis. Mixed features are extracted from the time series features by the TSF Mixer Block, and the MLP Block performs the final identification.

For clarity, the following description is given in reverse order, starting with the most minor parts.

Rotation Parameter Calculation Block ($n_{1\,\text{to}\,7}^{\text{bk}}$, $n_{1\,\text{to}\,7}^{\text{d}}$)



TSF Mixer Block ($n_{1\,\text{to}\,6}^{\text{bk}}$, $n_{1\,\text{to}\,6}^{\text{d}}$)

in

Flatten

MLP Block ($n_7^{\text{bk}}$, $n_7^{\text{d}}$)

FC (4)

tanh

out

Four parameters for Rodrigues' rotation formula

Fig. 6. Rotation Parameter Calculation Block

### 3.1 MLP Block

The structure of an MLP Block, which is the basic structural element used in various parts of rTsfNet, is shown in Figure 2. The number of stages of full-connection (FC) layers is $n$, and the number of kernels in the final FC layer is $n^{\text{base kernel}}$. Each FC layer has $n^{\text{bk}} \times 2^{n-i}$ kernels where $i$ shows its distance from the final layer. The output of the FC layers is activated with LeakyReLU after applying Layer Normalization. Next, a 50% dropout is applied.
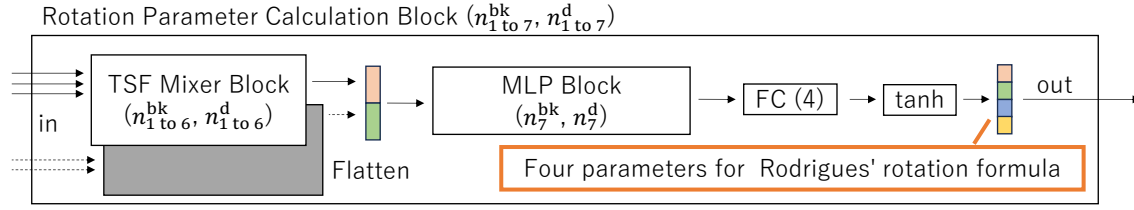
### 3.2 TSF Mixer sub-Block

The TSF Mixer sub-Block (Figure 3) is used in the TSF Mixer Block described in the next section. It has the structure of a TSF Mixer Block without the axis-wise and the channel-wise binary selections. It receives as input the data in which the TSFs for each axis are stored. First, an MLP Block is applied to extract the features within each axis. In this MLP Block, the weight is shared by every axis. The features derived from each axis by the MLP Block are serialized, and another MLP Block derives the final features.

### 3.3 TSF Mixer Block

The TSF Mixer Block (Figure 4) is a TSF Mixer sub-Block with axis-wise and channel-wise binary selections. The binary selection weights are calculated by the TSF Mixer sub-Blocks with input data. The calculated weights are applied to the mainstream after the axis-wise features are computed. Then we selectively set the values of the specific axes and channels to 0. This binary selection adaptively eliminates unnecessary axes from subsequent calculations among the axes augmented by the Multi-head 3D Rotation Block described below. Similarly, unnecessary TSF combinations are adaptively eliminated from subsequent calculations.

### 3.4 Multi-head 3D Rotation Block

The Multi-head 3D Rotation Block (Figure 5) calculates multiple 3D rotation parameters from the input 3D rotatable 3-axis sensor values and concatenates them after the rotations. First, we calculate and concatenate the L2 norm from the input. Then, we divided it into multiple block sets along the time series and extracted the TSFs from each. The extracted TSFs may be different for each block set. For example, an input with a data length of 128 is divided into a block set that consists of four blocks with a data length of 32 while deriving the mean and variance, and another block set that consists of one block with a data length of 128 while deriving the frequency feature. After the TSF extraction, tags are added to each axis. The tags consist of sensor location ID, sensor type, and axis type; they are simple integers, like 1, 2, and 3. All the derived TSFs are input into a single Rotation Parameter Calculation Block (Figure 6) to obtain Rodrigues' rotation formula parameters.

The Rotation Parameter Calculation Block first extracts the features for each input block set. Each segment in the block set is input to a TSF Mixer Block whose weight is shared within the block set. All the features extracted per block are serialized, and then the MLP Block derives the overall features. The four parameters required by Rodrigues' rotation formula, the XYZ elements of the rotation axis vectors, and a rotation angle are calculated by the FC layer with four kernels. The tanh was the most effective activation function of those parameters. The use of tanh is an empirical conclusion. Probably, the range limitation of tanh makes an effect.

Then the 3D rotation matrix is calculated from the derived Rodrigues' rotation formula parameters. The XYZ elements of the rotation axis vector are normalized to have an L2 norm of 1. Then the values of the 3-axis sensor, which is the input to Multi-head 3D Rotation, are rotated using the calculated 3D rotation matrix. This process is repeated for the number of heads, $n^{\text{head}}$. However, the second and subsequent Rodrigues' rotation formula parameters are obtained as the sum of those used before. This structure stabilizes the accuracy. The use of the sum is an empirical conclusion.

### 3.5 The main part of rTsfNet

Finally, the main part of rTsfNet, which uses these parts described above, is shown in Figure 1. In rTsfNet, the values of the 3-axis sensors that can rotate in 3D are treated separately from other values. Those former sensor values are subjected to multiple rotations in the Multi-head 3D Rotation Block. The same calculations as in the initial part of the Multi-head 3D Rotation Block are then performed to concatenate the L2 norm of the values of the 3-axis sensors that can rotate in 3D, the sensor values after the rotation, and the values of the sensors that cannot rotate in 3D. We then divide the data into several block sets along the time series, derive TSFs for each datum, and add a *tag* to each axis. Features are then extracted for each segment set. Each block in the block set is input to a TSF Mixer Block whose weight is shared within the block set. All the features extracted per block are serialized and the overall features are derived in the MLP Block. The identification results are obtained by FC layers with kernels of the number of classes to be identified and the Softmax activation function.

To summarize the structure of rTsfNet, rTsfNet has the following 29 number hyperparameters: $n^{\text{h}}$, $n^{\text{bk}}_{\text{1to14}}$, and $n^{\text{d}}_{\text{1to14}}$. It also needs block parameters such as block size, overlap size, and what TSFs will be extracted.

## 4 EVALUATION

Due to the verifiability issue mentioned in Section 2, we mainly performed comparisons using the UCI HAR Benchmark setup. The evaluation results of the iSPLInception-like benchmark setup, called the IMU-based HAR Benchmark setup, for PAMAP2 and Daphnet are also discussed. However, note that no strict comparison is possible due to the changes in the handling of the NaN values and the boundaries between trials. In addition, the evaluation results of the iSPLInception Benchmark setup for the UCI HAR and OPPORTUNITY datasets are also shown as reference values. The identification results for the newly defined IMU-based HAR Benchmark setup with the OPPORTUNITY dataset are also presented for future studies.

### 4.1 Meta settings of rTsfNet numeric hyperparameters

Due to exploration time issues, the following limitations were applied.

- $n^{\text{h}}_1 = 4$
- $n^{\text{bk}}_i = n^{\text{bk}}_j$, where $i = (2, 6)$, $j = (9, 13)$
- $n^{\text{d}}_i = n^{\text{d}}_j$, where $i = (2, 6)$, $j = (9, 13)$

Table 3. Selected time series features

| Description | Definition |
|---|---|
| Minimum | |
| Maximum | |
| Abs. energy | $\sum_{i=1}^{N} x_i^2$ |
| Abs. sum of changes | $\sum_{i=1}^{N} \lvert x_i - x_{i-1} \rvert$ |
| Mean change | $\frac{1}{N} \sum_{i=1}^{N} (x_i - x_{i-1})$ |
| Rooted mean squared | $\sqrt{\frac{1}{N} \sum_{i=1}^{N} x_i^2}$ |
| Count above start/end values of segments | |
| Number of crossings with 1st/3rd quartile. | |
| Mean of FFT amplitude ratio | |
| Skewness of FFT amplitude | |
| Mean of ac. with lags of multiple of N (only N=1 was used in this study) | |
| Kurtosis of ac. with lags of multiple of N (only N=1 was used in this study) | |

- $n_1^{\text{bk}} = n_3^{\text{bk}} = n_5^{\text{bk}} = n_8^{\text{bk}} = n_{10}^{\text{bk}} = n_{12}^{\text{bk}}$
- $n_1^{\text{d}} = n_3^{\text{d}} = n_5^{\text{d}} = n_8^{\text{d}} = n_{10}^{\text{d}} = n_{12}^{\text{d}}$

Different values were selected for the parameters $n_{4,11}^{\text{bk,d}}$ because the axis counts to which the calculated weights will be applied are significantly different. For a similar reason, different values are selected for the parameters $n_{7,14}^{\text{bk,d}}$ because the layer outputs are significantly different.

### 4.2 Block sizes

We used the following two block sizes with identical TSF extraction settings. One is a short block with a block size of about 0.5 seconds. There is no overlap between the blocks. The specific size depends on the sampling rate of the dataset. The other is a long block, where the entire segmentation is just one block. However, no long block was used in the iSPLInception Benchmark setup for OPPORTUNITY because it worsened the result due to the dirty segmentation issue.

### 4.3 Time series feature selection

This paper uses the time series features shown in Table 3 from the selection using genetic algorithms [7] and manual examination. The list of the time series features examined is shown in the Appendix. The selection was conducted for each dataset mentioned above, although identical features were selected for the final result. Features that were insufficiently defined in the table are described below.

*4.3.1 Mean of FFT amplitude ratio.* This ratio is the mean value of normalized frequency amplitude values derived by the FFT. The normalized range is 0 to 1.

*4.3.2 Mean/kurtosis of autocorrelation values with lags of multiples of N.* These features are the mean/kurtosis values of the autocorrelation values with lags of multiples of N, e.g., N = 2, 4, 6, 8 $\cdots$. The definition of the autocorrelation is $\frac{1}{(n-l)\sigma^2} \sum_{t=1}^{n-l}(X_t - \mu)(X_{t+l} - \mu)$. Only N=1 was used in this study. The maximum of the lag is limited to half of the segmentation, e.g., 16 if the segment consists of 32 samples.

### 4.4 Training proceedure

In this paper, we trained rTsfNet using the following settings.

- start learning late: 0.001
- max epochs: 350
- Reduce learning rate on training loss plateau: 20% decrease after 10 epochs in training loss plateau
- Early stop on validation loss plateau: 50 epochs
- Boot strap protection for early stop: 150 epochs

The final model was selected from a model with the best validation loss and the model of the final epoch because the model performance might be enhanced after the validation loss faces a plateau if the training loss was improved [11].

### 4.5 Evaluation with UCI HAR dataset benchmark setups

This section compares the performance of rTsfNet with other studies on the UCI HAR dataset benchmark setup.

*4.5.1 UCI HAR dataset.* The UCI HAR dataset [2] consists of two trials with 30 subjects each whose ages ranged from 19 and 48. A Samsung Galaxy S II smartphone was fixed on the left side of their waists for the first trial and on anywhere on their waists based on their own preferences for the second trial. The data from a three-axis accelerometer and a three-axis gyroscope were collected at 50 Hz with six activities: walking, walking upstairs, walking downstairs, sitting, standing, and lying. The data between each activity were removed from the dataset and then segmented into 128 samples with 50% overlap. The 70% of the subjects were selected for the train set and the rest for the test set.

The UCI HAR dataset satisfies all the benchmarking requirements.

*4.5.2 Performance comparison.* The rTsfNet's parameters for the UCI HAR are shown in Table 4, as selected by a genetic algorithm [7] and manual examination. However, not all the spaces have been fully explored due to computing times. A great possibility remains that better parameters exist.

As we have already seen, Table 1 summarizes the accuracy of the studies using the UCI HAR Benchmark setup, including the rTsfNet results. Table 5 shows a confusion matrix of the rTsfNet results.

rTsfNet achieved the highest accuracy, mf1, and wf1 results among the related studies. The second best result in Table 1, A. Ignatov [12], used data centering of the statistical features on the dataset, including the test set. Their result falls to the seventh place without the test set's information. The third place in Table 1, W. Jiang et al. [13], is an ensemble and boosting approach that combined of 2D FFT, CNN, and SVM. Thus, the highest performance among the related studies with end-to-end learning of DL is only in the fifth place, L. Lu et al. [20]. The rTsfNet performance is 1.09 points higher than that. It is a significant improvement.

*4.5.3 How well did the Multi-head 3D Rotation functon?* The fourth place position in Table 1 is held by rTsfNet without any Multi-head 3D Rotation. This pattern's parameters are shown in Table 4. This setup is identical as when TSFs are derived in advance and input to the DNN.

Its acc, mf1, and wf1 are 96.71, 0.9675, and 0.9670, respectively. This result shows that the effectiveness of combining the network structure of rTsfNet and the selected TSFs even without the Multi-head 3D Rotation. Moreover, the 97.76, 0.9779, and 0.9776 results for acc, mf1, and wf1 with Multi-head 3D Rotation show the overall improvements from this result. Thus, all of the network structures of rTsfNet, the selected TSFs, and Multi-head 3D Rotation are very effective in the IMU-based HAR.

Table 4. The rTsfNet parameters for each dataset on this study

| Parameter | UCI HAR | UCI HAR (no-mh-3D-rot.) | PAMAP2 | Daphnet | OPPORTUNITY (iSPL) | OPPORTUNITY |
|---|---|---|---|---|---|---|
| $n^{\mathrm{h}}$ | 4 | 4 | 4 | 4 | 4 | 4 |
| $n^{\mathrm{bk}}_{1,3,5,8,10,12}$ | 128 | 128 | 128 | 128 | 16 | 64 |
| $n^{\mathrm{d}}_{1,3,5,8,10,12}$ | 2 | 1 | 1 | 1 | 1 | 1 |
| $n^{\mathrm{bk}}_{2,9}$ | 128 | 128 | 32 | 16 | 64 | 64 |
| $n^{\mathrm{d}}_{2,9}$ | 3 | 4 | 3 | 1 | 1 | 3 |
| $n^{\mathrm{bk}}_{6,13}$ | 64 | 128 | 32 | 32 | 128 | 128 |
| $n^{\mathrm{d}}_{6,13}$ | 1 | 4 | 2 | 1 | 4 | 2 |
| $n^{\mathrm{bk}}_{4}$ | 128 | 16 | 32 | 16 | 16 | 64 |
| $n^{\mathrm{d}}_{4}$ | 4 | 3 | 3 | 3 | 3 | 1 |
| $n^{\mathrm{bk}}_{7}$ | 16 | 32 | 128 | 128 | 128 | 64 |
| $n^{\mathrm{d}}_{7}$ | 3 | 2 | 2 | 3 | 4 | 2 |
| $n^{\mathrm{bk}}_{11}$ | 16 | 64 | 128 | 64 | 16 | 16 |
| $n^{\mathrm{d}}_{11}$ | 4 | 4 | 4 | 2 | 1 | 3 |
| $n^{\mathrm{bk}}_{14}$ | 32 | 128 | 128 | 128 | 128 | 128 |
| $n^{\mathrm{d}}_{14}$ | 1 | 1 | 1 | 1 | 3 | 3 |
| Block size | 32, 128 | 32, 128 | 64, 256 | 32, 192 | 15 | 16, 32 |

Table 5. Confusion Matrix of rTsfNet for UCI HAR

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| A: Walking | 495 | 1 | 0 | 0 | 0 | 0 |
| B: Walking Upstairs | 4 | 466 | 1 | 0 | 0 | 0 |
| C: Walking Downstairs | 2 | 6 | 412 | 0 | 0 | 0 |
| D: Sitting | 0 | 0 | 0 | 462 | 29 | 0 |
| E: Standing | 0 | 0 | 0 | 23 | 509 | 0 |
| F: Laying | 0 | 0 | 0 | 0 | 0 | 537 |

## 4.6 Evaluation with PAMAP2, Daphnet, and Opportunity

To evaluate the method's generality of the method, rTsfNet was checked with PAMAP2, Daphnet, and OPPORTUNITY, all of which are more complex than UCI HAR. PAMAP2 is a dataset of basic behaviors under multi-sensor and complex conditions; OPPORTUNITY is a dataset of such daily behaviors as opening a door acquired with multi-sensors; Daphnet is a dataset targeting the freezing phenomenon of Parkinson's disease, which is more challenging to capture with accelerometers compared to basic behaviors.

The iSPLInception Benchmark setup clearly defines the segmentation and train/test set split for these datasets; however, its benchmark setup has a time warp issue in the segmentation due to the simple exclusion of samples containing NaN values and null labels and segmentation that does not consider the break between trials, as already

Table 6. IMU-based Benchmark setup details (+iSPL for OPPORTUNITY)

|  | UCI HAR | PAMAP2 | OPPORTUNITY | Daphnet |
|---|---|---|---|---|
| Sampling rate | 50 Hz | 100 Hz | 30 Hz | 64 Hz |
| Segment length | 128 | 256 | 32 (iSPL: 90) | 192 |
| Overlap | 50% | 50% | 50% | 50% |
| Training set | the others | the others | the others | the others |
| Validation set | 7, 22 | 5 | S1-ADL1, S3-ADL3,Drill, S4-ADL4 | S02R02, S03R03, S05R01 |
| Test set | 2,4,9,10,12,13,18,20,24 | 6 | S2-ADL2,Drill, S3-ADL1, S4-ADL5 | S02R01, S04R01, S05R02 |
| Split rule | subject-based | subject-based | trial-based | trial-based |

described in Subsection 2.3. Therefore, this study defines a new benchmark setup, called the IMU-based HAR Benchmark, with the following changes from the iSPLInception Benchmark setup.

- NaN values: linear interpolation up to 0.2 seconds
- Trial boundary: no segmentation over the boundary
- NULL label boundary: no segmentation over the boundary

In addition, with OPPORTUNITY, the segmentation length is shortened from 90 samples to 32 samples due to the activity length issue, also already described in Subsection 2.3.

Segmentation involving inter-user and inter-trial breaks is negligible and has little impact on all the datasets. The segmentation affected by simply excluding samples containing NaN values and null labels in PAMAP2 and Daphnet is also minimal and has little impact. Compared to these datasets, the OPPORTUNITY dataset significantly impacts this issue.

Thus, although the IMU-based HAR Benchmark and the iSPLInception Benchmark setups are not strictly comparable, the PAMAP2 and Daphnet data sets are less affected by the modifications, and so we use the values before and after the modifications for comparison and verification in this paper. On the other hand, since OPPORTUNITY is affected by a large number of modifications, we did not directly compare it. rTsfNet's results in the iSPLInception Benchmark setup are presented for comparison, but only as a reference value that should not be compared in future studies.

## 4.7 IMU-based HAR Benchmark

The IMU-based HAR Benchmark is an open-source platform defined in this study. It is forked from the iSPLInception Benchmark setup. Most setups are identical as the iSPLInception Benchmark. The following are the changes:

- NaN values: linear interpolation up to 0.2 seconds
- Trial boundary: no segmentation over the boundary
- NULL label boundary: no segmentation over the boundary
- Segmentation length of OPPORTUNITY: reduced from 90 to 32
- train/test set split of UCI HAR: returned to the original UCI HAR
- train/validation split of UCI HAR: newly defined.

The summary of each dataset's setup is shown in Table 6.[4]

---

[4]Although they were not used in this paper, we implemented supports for WISDM, RealWorld, the NULL label including segmentation for OPPORTUNITY, LOSO CV for datasets, Optuna based parameter optimization, etc., are implemented. Details are available: https://bit.ly/40b7R1C.

## 4.8 PAMAP2

The PAMAP2 dataset[22] consists of 11 basic and 7 optional activities with 9 subjects and three IMUs recorded at 100 Hz. IMU's locations are the hand, the chest, and the ankle. The IMUs recorded the readings of accelerometers, gyroscopes, magnetometers as well as the temperature and heart rates. The iSPLInception Benchmark and IMU-based HAR Benchmark setups use the data of accelerometers, gyroscopes, and magnetometers.

The following are the basic activities: lying, sitting, standing, walking, running, cycling, nordic walking, ascending stairs, descending stairs, vacuum cleaning, and ironing. The optional activities are watching TV, computer work, car driving, folding laundry, house cleaning, playing soccer, and rope jumping. The iSPLInception and IMU-based HAR Benchmark setups target the basic activities.

*4.8.1 Daphnet.* The Daphnet dataset[3] is a dataset to benchmark methods used for recognizing the Freezing of Gait (FOG) of Parkinson's disease (PD), which is a sudden and transient inability to walk. It places wearable acceleration sensors placed on the shank, the thigh, and the lower back. Almost 50% of patients with advanced Parkinson's disease face the FOG. It can cause of fails, interference with daily activities, and a significantly impaired quality of life.

This dataset consists of two targets: freeze and no freeze. The data were collected from 10 PD patients at a sampling rate of 64 Hz. Due to its disease incidence, this dataset is extremely class-imbalanced. Thus, this is a very difficult dataset to detect the target.

## 4.9 OPPORTUNITY

The Opportunity activity recognition dataset[23] targets 17 naturalistic activities in daily life, such as "Open Door", "Close Door", "Clean Table", "Drink from Cup" and so on. It was collected in a sensor-rich environment with 12 subjects. Due to frequent continuous NaN issues, the iSPLInception Benchmark and IMU-based HAR Benchmark setups used 7 wearable sensors of many: five XSens of the motion jacket, and two Sun SPOTs on the shoes. This dataset's sampling rate is 30 Hz. Four subjects recorded five trials for the activity set and one long remaining run to collect a large number of activity instances.

## 4.10 Performance evaluation

The rTsfNet parameters for PAMAP2, Daphnet, and OPPORTUNITY are shown in Table 4. They were selected by a genetic algorithm [7] and manual examination. However, not all the spaces have been completely explored due to computing times. Perpahs, better parameters will eventually be identified.

The comparison results are shown in Table 7. The result with the UCI HAR dataset on the iSPLInception Benchmark setup is also listed as a reference value. The rTsfNet showed the highest performance for all the datasets. The confusion matrixes of rTsfNet are shown in Tables 8, 9, 10, and 11.

The OPPORTUNITY result with the iSPLInception Benchmark setup is lower than OPPORTUNITY with the IMU-based HAR Benchmark setup. These two cannot be compared directly; however, the result suggests that the dirty segmentation of the iSPLInception Benchmark setup caused inaccurate classification.

## 4.11 Summary

rTsfNet showed the highest performance for all the datasets: UCI HAR, PAMAP2, Daphnet, and OPPORTUNITY, although each one has a different sensor setup and varying targets. This result means that rTsfNet's concept is suitable and has generality for IMU-based HAR.

Table 7. Comparision on the IMU-based HAR and iSPLInception Benchmark setups

| | PAMAP2 | | | Daphnet | | | UCI HAR (iSPL) | | |
| | acc | mf1 | wf1 | acc | mf1 | wf1 | acc | mf1 | wf1 |
|---|---|---|---|---|---|---|---|---|---|
| rTsfNet | 95.35 | 0.9353 | 0.9545 | 95.65 | 0.7051 | 0.9466 | 97.49** | 0.9749** | 0.9749** |
| iSPLInception[24] | 89.10 | 0.8786 | 0.8821 | 93.52 | 0.6533 | 0.9212 | 95.09 | 0.9499 | 0.9508 |
| CNN[35]* | 85.79 | 0.8424 | 0.8399 | 92.97 | 0.5455 | 0.9035 | 91.67 | 0.9160 | 0.9159 |
| CNN-LSTM[21]* | 88.37 | 0.8687 | 0.8720 | 92.97 | 0.5258 | 0.8993 | 94.48 | 0.9442 | 0.9441 |
| vLSTM[21]* | 85.53 | 0.8368 | 0.8487 | 93.22 | 0.5873 | 0.9106 | 90.80 | 0.9077 | 0.9082 |
| sLSTM[21]* | 88.42 | 0.8715 | 0.8866 | 87.65 | 0.5598 | 0.8797 | 91.82 | 0.9180 | 0.9174 |
| BiLSTM[33]* | 86.98 | 0.8597 | 0.8646 | 92.41 | 0.4803 | 0.8918 | 93.92 | 0.9383 | 0.9390 |

| | OPPORTUNITY (iSPL) | | | OPPORTUNITY | | |
| | acc | mf1 | wf1 | acc | mf1 | wf1 |
|---|---|---|---|---|---|---|
| rTsfNet | 91.76 | 0.8815 | 0.9177 | 94.65 | 0.9107 | 0.9462 |
| iSPLInception[24] | 88.14 | 0.8369 | 0.8811 | | | |
| CNN[35]* | 82.24 | 0.7384 | 0.8005 | | | |
| CNN-LSTM[21]* | 81.41 | 0.7375 | 0.8111 | | | |
| vLSTM[21]* | 76.79 | 0.6949 | 0.7676 | | | |
| sLSTM[21]* | 80.82 | 0.7194 | 0.8002 | | | |
| BiLSTM[33]* | 79.90 | 0.7297 | 0.7995 | | | |

\* No official implementation. It is based on the re-implementation of the iSPLIncepton Benchmark.
\*\* This model was trained with the original train/test split (meaning less training data).

Table 8. Confusion matrix of the PAMAP2

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A: Lying | 173 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 4 | 2 | 0 |
| B: Sitting | 0 | 46 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 5 | 0 |
| C: Standing | 0 | 1 | 34 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| D: Walking | 0 | 0 | 0 | 195 | 0 | 0 | 0 | 2 | 2 | 0 | 0 |
| E: Running | 0 | 0 | 0 | 2 | 166 | 0 | 0 | 3 | 6 | 0 | 0 |
| F: Cycling | 0 | 0 | 0 | 0 | 0 | 157 | 0 | 0 | 2 | 0 | 0 |
| G: Nordic walking | 0 | 1 | 0 | 1 | 0 | 0 | 204 | 1 | 0 | 0 | 0 |
| H: Ascending stairs | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 90 | 9 | 0 | 0 |
| I: Descending stairs | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 77 | 1 | 4 |
| J: Vacuum cleaning | 0 | 0 | 0 | 1 | 0 | 0 | 5 | 9 | 1 | 147 | 0 |
| K: Ironing | 0 | 0 | 2 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 289 |

Table 9. Confusion matrix of the Daphnet

| | A | B |
|---|---|---|
| A: No freeze | 2095 | 8 |
| B: Freeze | 89 | 37 |

As discussed in Subsection 4.5.3, although the combination of rTsfNet's network structure and the selected TSFs is effective even without Multi-head 3D Rotation, using it shows overall improvements from the results obtained by

Table 10.  Confusion matrix of the OPPORTUNITY (iSPL)

|                    | A  | B  | C  | D  | E  | F  | G  | H  | I  | J  | K  | L  | M  | N  | O  | P   | Q  |
|--------------------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|-----|----|
| A: Open Door 1     | 68 | 0  | 4  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 4   | 0  |
| B: Open Door 2     | 0  | 62 | 0  | 4  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| C: Close Door 1    | 2  | 0  | 72 | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| D: Close Door 2    | 1  | 2  | 0  | 61 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| E: Open Fridge     | 0  | 0  | 0  | 0  | 73 | 4  | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| F: Close Fridge    | 0  | 0  | 0  | 0  | 5  | 66 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| G: Open Dishwasher | 0  | 0  | 0  | 0  | 2  | 5  | 48 | 4  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| H: Close Dishwasher| 0  | 0  | 0  | 0  | 1  | 0  | 3  | 46 | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0  |
| I: Open Drawer 1   | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 1  | 43 | 5  | 2  | 0  | 0  | 0  | 0  | 0   | 0  |
| J: Close Drawer 1  | 0  | 0  | 0  | 0  | 0  | 2  | 0  | 0  | 1  | 19 | 2  | 1  | 0  | 0  | 0  | 0   | 3  |
| K: Open Drawer 2   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 1  | 1  | 29 | 2  | 3  | 0  | 0  | 0   | 0  |
| L: Close Drawer 2  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 2  | 2  | 21 | 1  | 0  | 0  | 0   | 0  |
| M: Open Drawer 3   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 2  | 3  | 37 | 1  | 0  | 0   | 0  |
| N: Close Drawer 3  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 2  | 3  | 5  | 38 | 0  | 0   | 0  |
| O: Clean Table     | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 97 | 0   | 1  |
| P: Drink from Cup  | 0  | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 1  | 260 | 0  |
| Q: Toggle Switch   | 0  | 2  | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 51 |

Table 11.  Confusion matrix of the OPPORTUNITY

|                    | A   | B   | C   | D   | E   | F   | G   | H  | I  | J  | K  | L  | M  | N  | O   | P   | Q  |
|--------------------|-----|-----|-----|-----|-----|-----|-----|----|----|----|----|----|----|----|-----|-----|----|
| A: Open Door 1     | 146 | 0   | 15  | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 4   | 0  |
| B: Open Door 2     | 0   | 143 | 0   | 2   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| C: Close Door 1    | 2   | 0   | 166 | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| D: Close Door 2    | 0   | 2   | 0   | 151 | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 1  |
| E: Open Fridge     | 0   | 0   | 0   | 0   | 156 | 6   | 0   | 0  | 2  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| F: Close Fridge    | 0   | 0   | 0   | 0   | 7   | 116 | 0   | 2  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| G: Open Dishwasher | 0   | 0   | 0   | 0   | 0   | 0   | 105 | 8  | 13 | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| H: Close Dishwasher| 0   | 0   | 0   | 0   | 0   | 0   | 11  | 96 | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 0  |
| I: Open Drawer 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 76 | 0  | 7  | 0  | 0  | 0  | 0   | 0   | 5  |
| J: Close Drawer 1  | 0   | 0   | 0   | 0   | 0   | 0   | 1   | 2  | 5  | 27 | 0  | 5  | 0  | 0  | 0   | 0   | 3  |
| K: Open Drawer 2   | 0   | 0   | 0   | 0   | 0   | 0   | 1   | 1  | 5  | 0  | 48 | 0  | 0  | 0  | 0   | 0   | 0  |
| L: Close Drawer 2  | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 0  | 6  | 1  | 29 | 0  | 1  | 0   | 0   | 0  |
| M: Open Drawer 3   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 3  | 0  | 65 | 3  | 0   | 0   | 0  |
| N: Close Drawer 3  | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 2  | 1  | 91 | 0   | 0   | 0  |
| O: Clean Table     | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 235 | 0   | 0  |
| P: Drink from Cup  | 1   | 0   | 0   | 0   | 0   | 0   | 2   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 2   | 656 | 0  |
| Q: Toggle Switch   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0   | 0   | 99 |

rTsfNet without it. All the network structures of rTsfNet, the selected TSFs, and the Multi-head 3D Rotation are very effective in IMU-based HAR domain.

## 5   POTENTIAL

Next, we described rTsfNet's potential.

**Cooperation with extensions**  rTsfNet can probably be improved in combination with such extension structures as Residual, SE, Attention, LSTM, ensemble, and so on. We propose rTsfNet in this paper as a basic structure for other networks like CNN.

**As feature extractor**  Since rTsfNet can be used as a feature extractor. Therefore, for example, it is likely possible that VAE-like networks can improve their performance with the use of rTsfNet as their encoder.

**TSF implementation**  The TSFs described in this paper are only a few proposed by many researchers. rTsfNet's performance can be improved if more complex, suitable, or lightweight TSFs are implemented into the network.

**TSF and parameter selection**  The TSFs and parameters of rTsfNet described in this paper were selected by a genetic algorithm [7] and manual examination. However, not every space has been fully explored due to the required computing times. Perhaps even better parameters can be identified. In addition, several parameters in this study have the same values to reduce exploration times, as described in subsection 4.1. If such limitations were removed, their performance would be increased.

## 6    CONCLUSION

As a new DNN model for IMU-based HAR, this paper presented rTsfNet, a DNN model with Multi-head 3D Rotation and Time Series Feature Extraction. rTsfNet automatically selects 3D bases from which features should be derived by extracting 3D rotation parameters within the DNN. Time series features (TSFs) (which embody the wisdom of many researchers) are derived for achieving HAR using MLP.

Although our model does not use CNN, it achieved higher accuracy than the existing models under multiple datasets, which target different activities, under well-managed benchmark conditions. rTsfNet's concept is suitable and has generality for IMU-based HAR.

As discussed in Subsection 4.5.3, although the combination of the network structure of rTsfNet and the selected TSFs is effective even without the use of the multi-head 3D rotation, using it shows overall improvements from the results obtained by rTsfNet without it. It means that all of the network structures of rTsfNet, the selected TSFs, and Multi-head 3D rotation are very effective in IMU-based HAR.

As an additional contribution, this study newly defined an IMU-based HAR Benchmark setup for these datasets and created direct comparability for studies that use the benchmark setup.

The rTsfNet source code and trained models are available at https://bit.ly/40b7R1C. The IMU-based HAR Benchmark system is available at https://bit.ly/45OZ1aT.

## ACKNOWLEDGMENTS

## REFERENCES

[1]  Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge L. Reyes-Ortiz. 2012.  Human Activity Recognition on Smartphones Using a Multiclass Hardware-Friendly Support Vector Machine. In *Ambient Assisted Living and Home Care*, José Bravo, Ramón Hervás, and Marcela Rodríguez (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 216–223.

[2] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, Jorge Luis Reyes-Ortiz, et al. 2013. A public domain dataset for human activity recognition using smartphones.. In *Esann*, Vol. 3. 3.

[3] Marc Bachlin, Meir Plotnik, Daniel Roggen, Inbal Maidan, Jeffrey M. Hausdorff, Nir Giladi, and Gerhard Troster. 2010. Wearable Assistant for Parkinson's Disease Patients With the Freezing of Gait Symptom. *IEEE Transactions on Information Technology in Biomedicine* 14, 2 (2010), 436–446. https://doi.org/10.1109/TITB.2009.2036165

[4] Gustavo EAPA Batista, Eamonn J Keogh, Oben Moses Tataw, and Vinicius MA De Souza. 2014. CID: an efficient complexity-invariant distance for time series. *Data Mining and Knowledge Discovery* 28 (2014), 634–669.

[5] Valentina Bianchi, Marco Bassoli, Gianfranco Lombardo, Paolo Fornacciari, Monica Mordonini, and Ilaria De Munari. 2019. IoT Wearable Sensor and Deep Learning: An Integrated Approach for Personalized Human Activity Recognition in a Smart Home Environment. *IEEE Internet of Things Journal* 6, 5 (2019), 8553–8562. https://doi.org/10.1109/JIOT.2019.2920283

[6] Abdelghani Dahou, Mohammed A.A. Al-qaness, Mohamed Abd Elaziz, and Ahmed Helmi. 2022. Human activity recognition in IoHT applications using Arithmetic Optimization Algorithm and deep learning. *Measurement* 199 (2022), 111445. https://doi.org/10.1016/j.measurement.2022.111445

[7] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. 2002. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation* 6, 2 (2002), 182–197. https://doi.org/10.1109/4235.996017

[8] Yilin Dong, Xinde Li, Jean Dezert, Mohammad Omar Khyam, Md Noor-A-Rahim, and Shuzhi Sam Ge. 2020. DSmT-Based Fusion Strategy for Human Activity Recognition in Body Sensor Networks. *IEEE Transactions on Industrial Informatics* (Feb. 2020). https://doi.org/10.1109/TII.2020.2976812

[9] Harish Haresamudram, Irfan Essa, and Thomas Plötz. 2022. Assessing the State of Self-Supervised Human Activity Recognition Using Wearables. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3, Article 116 (sep 2022), 47 pages. https://doi.org/10.1145/3550299

[10] Jun He, Qian Zhang, Liqun Wang, and Ling Pei. 2019. Weakly Supervised Human Activity Recognition From Wearable Sensors by Recurrent Attention Learning. *IEEE Sensors Journal* 19, 6 (2019), 2287–2297. https://doi.org/10.1109/JSEN.2018.2885796

[11] Elad Hoffer, Itay Hubara, and Daniel Soudry. 2017. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/a5e0ff62be0b08456fc7f1e88812af3d-Paper.pdf

[12] Andrey Ignatov. 2018. Real-time human activity recognition from accelerometer data using Convolutional Neural Networks. *Applied Soft Computing* 62 (2018), 915–922. https://doi.org/10.1016/j.asoc.2017.09.027

[13] Wenchao Jiang and Zhaozheng Yin. 2015. Human Activity Recognition Using Wearable Sensors by Deep Convolutional Neural Networks. In *Proceedings of the 23rd ACM International Conference on Multimedia* (Brisbane, Australia) (*MM '15*). Association for Computing Machinery, New York, NY, USA, 1307–1310. https://doi.org/10.1145/2733373.2806333

[14] Hua Kang, Qianyi Huang, and Qian Zhang. 2022. Augmented Adversarial Learning for Human Activity Recognition with Partial Sensor Sets. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3, Article 122 (sep 2022), 30 pages. https://doi.org/10.1145/3550285

[15] Zanobya N. Khan and Jamil Ahmad. 2021. Attention induced multi-head convolutional neural network for human activity recognition. *Applied Soft Computing* 110 (2021), 107671. https://doi.org/10.1016/j.asoc.2021.107671

[16] Yong-Joong Kim, Bong-Nam Kang, and Daijin Kim. 2015. Hidden Markov Model Ensemble for Activity Recognition Using Tri-Axis Accelerometer. In *2015 IEEE International Conference on Systems, Man, and Cybernetics*. 3036–3041. https://doi.org/10.1109/SMC.2015.528

[17] Jennifer R. Kwapisz, Gary M. Weiss, and Samuel A. Moore. 2011. Activity Recognition Using Cell Phone Accelerometers. *SIGKDD Explor. Newsl.* 12, 2 (mar 2011), 74–82. https://doi.org/10.1145/1964897.1964918

[18] Haobo Li, Aman Shrestha, Hadi Heidari, Julien Le Kernec, and Francesco Fioranelli. 2020. Bi-LSTM Network for Multimodal Continuous Human Activity Recognition and Fall Detection. *IEEE Sensors Journal* 20, 3 (2020), 1191–1201. https://doi.org/10.1109/JSEN.2019.2946095

[19] Yongmou Li, Dianxi Shi, Bo Ding, and Dongbo Liu. 2014. Unsupervised Feature Learning for Human Activity Recognition Using Smartphone Sensors. In *Mining Intelligence and Knowledge Exploration*, Rajendra Prasath, Philip O'Reilly, and T. Kathirvalavakumar (Eds.). Springer International Publishing, Cham, 99–107.

[20] Limeng Lu, Chuanlin Zhang, Kai Cao, Tao Deng, and Qianqian Yang. 2022. A Multichannel CNN-GRU Model for Human Activity Recognition. *IEEE Access* 10 (2022), 66797–66810. https://doi.org/10.1109/ACCESS.2022.3185112

[21] Ronald Mutegeki and Dong Seog Han. 2020. A CNN-LSTM Approach to Human Activity Recognition. In *2020 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*. 362–366. https://doi.org/10.1109/ICAIIC48513.2020.9065078

[22] Attila Reiss and Didier Stricker. 2012. Introducing a New Benchmarked Dataset for Activity Monitoring. In *2012 16th International Symposium on Wearable Computers*. 108–109. https://doi.org/10.1109/ISWC.2012.13

[23] Daniel Roggen, Alberto Calatroni, Mirco Rossi, Thomas Holleczek, Kilian Förster, Gerhard Tröster, Paul Lukowicz, David Bannach, Gerald Pirkl, Alois Ferscha, Jakob Doppler, Clemens Holzmann, Marc Kurz, Gerald Holl, Ricardo Chavarriaga, Hesam Sagha, Hamidreza Bayati, Marco Creatura, and José del R. Millàn. 2010. Collecting complex activity datasets in highly rich networked sensor environments. In *2010 Seventh International Conference on Networked Sensing Systems (INSS)*. 233–240. https://doi.org/10.1109/INSS.2010.5573462

[24] Mutegeki Ronald, Alwin Poulose, and Dong Seog Han. 2021. iSPLInception: An Inception-ResNet Deep Learning Architecture for Human Activity Recognition. *IEEE Access* 9 (2021), 68985–69001. https://doi.org/10.1109/ACCESS.2021.3078184

[25] Charissa Ann Ronao and Sung-Bae Cho. 2014. Human activity recognition using smartphone sensors with two-stage continuous hidden Markov models. In *2014 10th International Conference on Natural Computation (ICNC)*. 681–686. https://doi.org/10.1109/ICNC.2014.6975918

[26] Charissa Ann Ronao and Sung-Bae Cho. 2015. Deep Convolutional Neural Networks for Human Activity Recognition with Smartphone Sensors. In *Neural Information Processing*, Sabri Arik, Tingwen Huang, Weng Kin Lai, and Qingshan Liu (Eds.). Springer International Publishing, Cham, 46–53.

[27] Charissa Ann Ronao and Sung-Bae Cho. 2016. Human activity recognition with smartphone sensors using deep learning neural networks. *Expert Systems with Applications* 59 (2016), 235–244. https://doi.org/10.1016/j.eswa.2016.04.032

[28] Charissa Ann Ronao and Sung-Bae Cho. 2017. Recognizing human activities from smartphone sensors using hierarchical continuous hidden Markov models. *International Journal of Distributed Sensor Networks* 13, 1 (2017), 1550147716683687. https://doi.org/10.1177/1550147716683687 arXiv:https://doi.org/10.1177/1550147716683687

[29] Skyler Seto, Wenyu Zhang, and Yichen Zhou. 2015. Multivariate Time Series Classification Using Dynamic Time Warping Template Selection for Human Activity Recognition. In *2015 IEEE Symposium Series on Computational Intelligence*. 1399–1406. https://doi.org/10.1109/SSCI.2015.199

[30] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going Deeper With Convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[31] Timo Sztyler and Heiner Stuckenschmidt. 2016. On-body localization of wearable devices: An investigation of position-aware activity recognition. In *2016 IEEE International Conference on Pervasive Computing and Communications (PerCom)*. 1–9. https://doi.org/10.1109/PERCOM.2016.7456521

[32] Kei Tanigaki, Tze Chuin Teoh, Naoya Yoshimura, Takuya Maekawa, and Takahiro Hara. 2022. Predicting Performance Improvement of Human Activity Recognition Model by Additional Data Collection. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3, Article 142 (sep 2022), 33 pages. https://doi.org/10.1145/3550319

[33] Nguyen Thi Hoai Thu and Dong Seog Han. 2020. Utilization of Postural Transitions in Sensor-based Human Activity Recognition. In *2020 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*. 177–181. https://doi.org/10.1109/ICAIIC48513.2020.9065250

[34] Nilay Tufek, Murat Yalcin, Mucahit Altintas, Fatma Kalaoglu, Yi Li, and Senem Kursun Bahadir. 2020. Human Action Recognition Using Deep Learning Methods on Limited Sensory Data. *IEEE Sensors Journal* 20, 6 (2020), 3101–3112. https://doi.org/10.1109/JSEN.2019.2956901

[35] D. van Kuppevelt, C. Meijer, F. Huber, A. van der Ploeg, S. Georgievska, and V.T. van Hees. 2020. Mcfly: Automated deep learning on time series. *SoftwareX* 12 (2020), 100548. https://doi.org/10.1016/j.softx.2020.100548

[36] Kun Wang, Jun He, and Lei Zhang. 2019. Attention-Based Convolutional Neural Network for Weakly Labeled Human Activities' Recognition With Wearable Sensors. *IEEE Sensors Journal* 19, 17 (2019), 7598–7604. https://doi.org/10.1109/JSEN.2019.2917225

[37] Kun Xia, Jianguang Huang, and Hanyu Wang. 2020. LSTM-CNN Architecture for Human Activity Recognition. *IEEE Access* 8 (2020), 56855–56866. https://doi.org/10.1109/ACCESS.2020.2982225

[38] Cheng Xu, Duo Chai, Jie He, Xiaotong Zhang, and Shihong Duan. 2019. InnoHAR: A Deep Neural Network for Complex Human Activity Recognition. *IEEE Access* 7 (2019), 9893–9902. https://doi.org/10.1109/ACCESS.2018.2890675

[39] Yu Zhao, Rennong Yang, Guillaume Chevalier, Ximeng Xu, and Zhenxing Zhang. 2018. Deep residual bidir-LSTM for human activity recognition using wearable sensors. *Mathematical Problems in Engineering* 2018 (2018), 1–13.

## A TIME SERIES FEATURES CONSIDERED IN THIS PAPER

The time series features that are used for the genetic algorithm-based selections are shown in Tables 12 and 13. We omitted the definitions of well-known features.

Table 12. Considered time series features (1/2)

| No. | description | |
|---|---|---|
| 1 | mean | |
| 2 | minimum | |
| 3 | maximum | |
| 4 | quantiles | the 1st, 2nd (median), and 3rd quartile. |
| 4 | time based quantiles | the values of 25%, 50%, and 75% point along time series. |
| 5 | skewness | |
| 6 | kurtosis | |
| 7 | variance | |
| 8 | standard deviation | |
| 9 | rooted mean squared | $\sqrt{\frac{1}{N}\sum_{i=1}^{N} x_i^2}$ |
| 10 | mean change | $\frac{1}{N}\sum_{i=1}^{N}(x_i - x_{i-1})$ |

Table 13. Considered time series features (2/2)

| No. | description | |
| --- | --- | --- |
| 11 | sum of change | $\sum_{i=1}^{N}(x_i - x_{i-1})$ |
| 12 | mean abs. change | $\frac{1}{N}\sum_{i=1}^{N}|x_i - x_{i-1}|$ |
| 13 | abs. energy | $\sum_{i=1}^{N} x_i^2$ |
| 14 | abs. sum of changes | $\sum_{i=1}^{N}|x_i - x_{i-1}|$ |
| 15 | abs. max | $\max_{1 \le i \le N}|x_i|$ |
| 16 | CID [4] | $\sqrt{\sum_{i=1}^{N-1}(x_i - x_{i-1})^2}$ |
| 17 | count above zero | |
| 18 | count above segment's mean | |
| 19 | count above segment's start value | |
| 20 | count above value of 25% position along segment's time series | |
| 21 | count above value of 50% position along segment's time series | |
| 22 | count above value of 75% position along segment's time series | |
| 23 | count above segment's end value | |
| 24 | number of crossings with zero | |
| 25 | number of crossings with segment's mean | |
| 26 | number of crossings with segment's 1st quantile | |
| 27 | number of crossings with segment's 2nd quantile | |
| 28 | number of crossings with segment's 3rt quantile | |
| 29 | number of crossings with segment's start value | |
| 30 | number of crossings with segment's value of 25% position along segment's time series | |
| 31 | number of crossings with value of 50% position along segment's time series | |
| 32 | number of crossings with value of 75% position along segment's time series | |
| 33 | number of crossings with segment's end value | |
| 34 | FFT amplitude | |
| 35 | FFT amplitude ratio | |
| 36 | mean of FFT amplitude | |
| 37 | variance of FFT amplitude | |
| 38 | skewness of FFT amplitude | |
| 39 | kurtosis of FFT amplitude | |
| 40 | mean of FFT amplitude ratio | |
| 41 | variance of FFT amplitude ratio | |
| 42 | skewness of FFT amplitude ratio | |
| 43 | kurtosis of FFT amplitude ratio | |
| 44 | FFT angle | |
| 45 | autocorrelation | $\frac{1}{(n-l)\sigma^2}\sum_{t=1}^{n-l}(X_t - \mu)(X_{t+l} - \mu)$ |
| 46 | mean of autocorrelation with lags of multiple of N | |
| 47 | variance of autocorrelation with lags of multiple of N | |
| 48 | skewness of autocorrelation with lags of multiple of N | |
| 49 | kurtosis of autocorrelation with lags of multiple of N | |