
Scale-teaching: Robust Multi-scale Training for Time Series Classification with Noisy Labels

Zhen Liu

South China University of Technology
Guangzhou, China
cszhenliu@mail.scut.edu.cn

Peitian Ma

South China University of Technology
Guangzhou, China
ma_scuter@163.com

Dongliang Chen

South China University of Technology
Guangzhou, China
ytucdl@foxmail.com

Wenbin Pei

Dalian University of Technology
Dalian, China
peiwenbin@dlut.edu.cn

Qianli Ma*

South China University of Technology
Guangzhou, China
qianlima@scut.edu.cn

Abstract

Deep Neural Networks (DNNs) have been criticized because they easily overfit noisy (incorrect) labels. To improve the robustness of DNNs, existing methods for image data regard samples with small training losses as correctly labeled data (small-loss criterion). Nevertheless, time series' discriminative patterns are easily distorted by external noises (i.e., frequency perturbations) during the recording process. This results in training losses of some time series samples that do not meet the small-loss criterion. Therefore, this paper proposes a deep learning paradigm called *Scale-teaching* to cope with time series noisy labels. Specifically, we design a fine-to-coarse cross-scale fusion mechanism for learning discriminative patterns by utilizing time series at different scales to train multiple DNNs simultaneously. Meanwhile, each network is trained in a cross-teaching manner by using complementary information from different scales to select small-loss samples as clean labels. For unselected large-loss samples, we introduce multi-scale embedding graph learning via label propagation to correct their labels by using selected clean samples. Experiments on multiple benchmark time series datasets demonstrate the superiority of the proposed Scale-teaching paradigm over state-of-the-art methods in terms of effectiveness and robustness.

1 Introduction

Time series classification has recently received much attention in deep learning [1, 2]. Essentially, the success of Deep Neural Networks (DNNs) is driven by a large amount of well-labeled data. However, human errors [3] and sensor failures [4] produce noisy (incorrect) labels in time series datasets. For example, in electrocardiogram diagnosis [5], physicians with different experiences tend to make inconsistent category judgments. In recent studies [6, 7], DNNs have shown their powerful learning ability, which, however, makes it relatively easier to overfit noisy labels and inevitably degenerate

*Qianli Ma is the corresponding author.

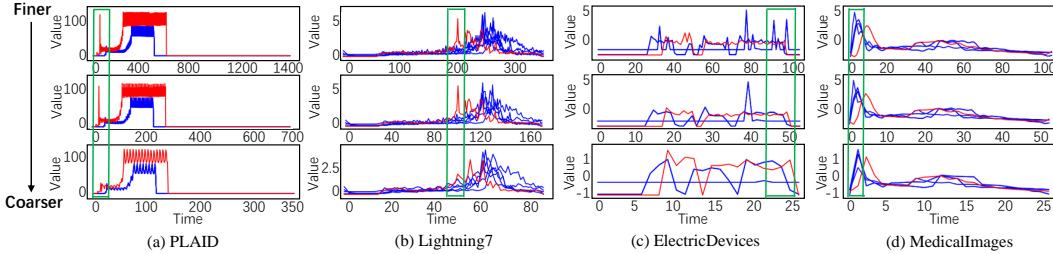


Figure 1: Illustration of time series samples *from the same category* at different time scales. Among all samples in the same category, **red** indicates the one with the largest variance, and **blue** indicates a few samples with the smallest variance.

the robustness of models. Moreover, time series data has complex temporal dynamics that make it challenging to manually correct noisy labels [8].

To cope with noisy labels, existing studies on label-noise learning [9, 10] use the memory effect of DNNs to select samples with small losses for training. DNNs memorize the data with clean labels first, and then those with noisy labels in classification training (small-loss criterion) [11]. It is worth noting that the small-loss criterion is not affected by the choice of training optimizations and network structures [12], and is widely utilized for label-noise learning in computer vision [13, 14]. However, the small loss criterion cannot always be applied to time series because the discriminative patterns of time series data are easily distorted by external noises [15, 16]. For example, in a smart grid, distortions may occur due to sampling frequency perturbations, imprecise sensors, or random differences in energy consumption [17]. Such distortions can make it difficult for DNNs to learn the appropriate discriminant patterns of time series, resulting in large training losses for some clean labeled samples. In addition, the small-loss criterion only utilizes the data’s label information and does not consider the inherent properties of time series features (i.e., multi-scale information).

Multi-scale properties are crucial in time series classification tasks. In recent years, multi-scale convolution [16], dynamic skip connections [18, 19] and adaptive convolution kernel size [20] have been utilized to learn discriminative patterns of time series. Furthermore, according to related studies [2, 20, 21], the selection of appropriate time scales for time series data can facilitate DNNs to learn class-characteristic patterns. With correct labels, the above studies indicate that the multi-scale properties of time series data can help DNNs learn appropriate discriminative patterns for mitigating the negative effects of time series recording noises. Nevertheless, it remains an open challenge as to how the multi-scale properties of time series can be used for label-noise learning.

To this end, we propose a deep learning paradigm, named Scale-teaching, for time-series classification with noisy labels. In particular, we design a fine-to-coarse cross-scale fusion mechanism for obtaining robust time series embeddings in the presence of noisy labels. We select four time series datasets from the UCR archive [22] to explain our motivation. As shown in Figure 1, in the single scale case (top row), the red and blue samples from the same category have large differences in certain local regions (the green rectangle in Figure 1). By downsampling the time series from fine to coarse, some local regions between the red and blue samples did become similar. Meanwhile, existing studies [12, 23] show that multiple DNNs with random initialization have classification divergence for noisy labeled samples, but are consistent for clean labeled samples. The above findings inspire us to utilize multiple DNNs to combine robust embeddings at different scales to deal with noisy labels. Nonetheless, the coarse scale discards many local regions in the fine scale (as in Figure 1 (c)), which may degenerate the classification performance. Hence, we propose the Scale-teaching paradigm, which can better preserve the local discriminative patterns of fine scale while dealing with distortions.

More specifically, the proposed Scale-teaching paradigm performs the cross-scale embedding fusion in the finer-to-coarser direction by utilizing time series at different scales to train multiple DNNs simultaneously. The cross-scale embedding fusion exploits complementary information from different scales to learn discriminative patterns. This enables the learned embeddings to be more robust to distortions and noisy labels. During training, clean labels are selected through cross-teaching on those networks with the learned embeddings. The small-loss samples in training are used as (clean) labeled data, and the unselected large-loss samples are used as (noisy) unlabeled data. Moreover,

multi-scale embedding graph learning is introduced to establish relationships between labeled and unlabeled samples for noisy label correction. Based on the multi-scale embedding graph, the label propagation theory [24] is employed to correct noisy labels. This drives the model to better fit time series category distribution. The contributions are summarized as follows:

- We propose a deep learning paradigm, called Scale-teaching, for time-series label-noise learning. In particular, a cross-scale fusion mechanism is designed to help the model select more reliable clean labels by exploiting complementary information from different scales.
- We further introduce multi-scale embedding graph learning for noisy label correction using the selected clean labels based on the label propagation theory. Unlike conventional image label-noise learning methods focused on sample loss levels, our approach uses well-learned multi-scale time series embeddings for noise label correction at sample feature levels.
- Extensive experiments on multiple benchmark time series datasets show that the proposed Scale-teaching paradigm achieves a state-of-the-art classification performance. In addition, multi-scale analyses and ablation studies indicate that the use of multi-scale information can effectively improve the robustness of Scale-teaching against noisy labels.

2 Related Work

Label-noise Learning. Existing label-noise learning studies focus mainly on image data [10]. These studies can be broadly classified into three categories: (1) designing noise-robust objective functions [25, 26] or regularization strategies [27, 28]; (2) detecting and correcting noisy labels [13, 29, 30]; (3) transition-matrix-based [31, 32] and semi-supervised-based [14, 33] methods. In contrast to the methodologies in the first and third categories, approaches categorized under the second category have received considerable attention in recent years [7, 34]. Methods of the second category can be further divided into sample selection and label correction. The common methods of sample selection are the Co-teaching family [12, 13, 23] and FINE [35]. Label correction [36, 37] attempts to correct noisy labels by either using prediction results of classifiers or pseudo-labeling techniques. Recently, SREA [4] utilizes pseudo-labels generated based on a clustering task to correct time-series noisy labels. Although the above methods can improve the robustness of DNNs, how the multi-scale properties of time series are exploited for label-noise learning has not been explored.

Multi-scale Time Series Modeling. In recent years, multi-scale properties have gradually gained attention in various time series downstream tasks [18, 38], such as time series classification, prediction, and anomaly detection [39]. For example, Cui et al. [16] employ multiple convolutional network channels of different scales to learn temporal patterns that facilitate time series classification. Chen et al. [19] design a time-aware multi-scale RNN model for human action prediction. Wang et al. [40] introduce a multi-scale one-class RNN for time series anomaly detection. Also, recent studies [41, 42, 43, 44] indicate that multi-scale properties can effectively improve the performance of long-term time series prediction. Unlike prior work, we utilize multiple DNNs with identical architectures to separately capture discriminative temporal patterns across various scales. This enables us to acquire robust embeddings for handling noisy labels via a cross-scale fusion strategy.

Label Propagation. Label propagation (LP) is a graph-based inductive inference method [24, 45] that can propagate pseudo-labels to unlabeled graph nodes using labeled graph nodes. Since LP can utilize the feature information of data to obtain pseudo-labels of unlabeled samples, related works employ LP in few-shot learning [46] and semi-supervised learning [47, 48]. Generally speaking, DNNs have the powerful capability for feature extraction, and the learned embeddings tend to be similar within classes and different between classes. Each sample contains feature and label information. Intuitively, the embeddings of samples with noisy labels obtained by DNNs closely align with the true class distribution when the DNNs do not fit noisy labels in the early training stages. Naturally, we create a nearest-neighbor graph based on well-learned multi-scale time series embeddings at the feature level. Subsequently, we employ LP theory to correct the labels of unselected noisy samples using the labels of clean samples chosen by the DNNs. This approach leverages robust multi-scale embeddings to address the issue of noisy labels.

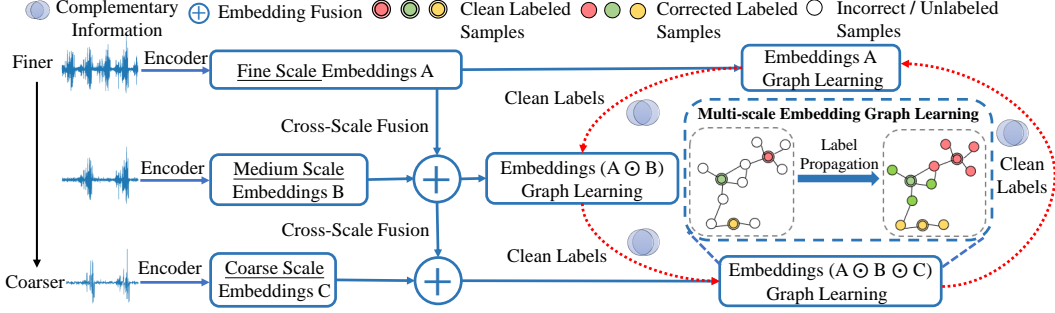


Figure 2: The Scale-teaching paradigm’s general architecture comprises two core processes: (i) clean label selection and (ii) noisy label correction. In the clean label selection phase, networks A, B, and C engage in cross-scale fusion, moving from fine to coarse (A→B, B→C). They employ clean labels acquired through cross-teaching (A→B, B→C, C→A) to guide their respective classification training. In the noisy label correction phase, pseudo labels derived from multi-scale embeddings graph learning are employed as corrected labels for time series not selected as clean labeled samples.

3 Proposed Approach

3.1 Problem Definition

Given a noisy labeled time series dataset $\mathcal{D} = \{(\mathcal{X}_i, \hat{y}_i)\}_{i=1}^N$, it contains N time series, where $\mathcal{X}_i \in R^{L \times T}$, L denotes the number of variables, and T is the length of variable. $\hat{y}_i \in \{1, \dots, C\}$ is the observed label of \mathcal{X}_i with η probability of being a noisy label. Our goal is to enable the DNNs trained on the noisy labeled training dataset \mathcal{D}_{train} to correctly predict the ground-truth labels of the given time series in the test set. Specifically, the problem to be addressed in this paper consists of two steps. The first is to select clean labeled time series from \mathcal{D}_{train} , and the second is to perform noisy label correction for time series in \mathcal{D}_{train} that have not been selected as clean labels.

3.2 Model Architecture

The overall architecture of Scale-teaching is shown in Figure 2. While this figure illustrates Scale-teaching with input time series at three scales, it can be extended to models with more scales, exceeding three. We utilize a consistent structural encoder to learn embeddings for each input scale sequence. Each encoder undergoes training at two levels: embedding learning for clean sample selection at the feature level and label correction with the multi-scale embeddings. For embedding learning, we propose a cross-scale fusion (Section 3.3) mechanism from fine to coarse to obtain robust embeddings. This approach enables the selection of more dependable clean labels through the small-loss criterion. Specifically, embeddings (A ⊗ B ⊗ C) encompass multi-scale information from fine, medium, and coarse scale sequences derived from the same time series. Regarding noisy label correction, we introduce multi-scale embedding graph learning (Section 3.4) based on label propagation, utilizing the selected clean samples to correct the labels of unselected large-loss samples.

3.3 Cross-scale Fusion for Clean Label Selection

After downsampling the original time series at different scales, it eliminates some of the differences in local regions between samples of the same category (as in Figure 1). However, the downsampled sequences (i.e., coarse scale) discard many local regions of the original time series. This tends to degrade the model’s classification performance if the downsampled sequences are used directly for classification (please refer to Table 2 in the Experiments section). Meanwhile, existing studies [12, 23] on label-noise learning show that DNNs with different random initializations have high consistency in classification results for clean labeled samples in the early training period, while there is disagreement in the classification of noisy labeled samples. Based on the above findings, we utilize multiple DNNs (or encoders) with different random initializations to learn embeddings of different downsampled scale sequences separately, and perform cross-scale fusion. On the one hand, we exploit complementary

information between adjacent scale embeddings to promote learned embeddings to be more robust for classification. On the other hand, we leverage the divergence in the classification of noisy labeled samples by different DNNs to mitigate the negative impact of noisy labels in training. In this way, we can utilize the cross-scale fusion embeddings for classification, thus better using the small loss criterion [11, 29] for clean label selection. Specifically, downsampling is employed to generate different scale sequences from the same time series. Given a time series $\mathcal{X}_i = \{x_1, x_2, \dots, x_T\}$, supposing the downsampling ratio is k . Then, we only keep data points in \mathcal{X}_i as follows:

$$\mathcal{X}_i^k = \{x_{k*j}\}, j = 1, 2, \dots, \frac{T}{k}, \quad (1)$$

where $k \in [1, T/2]$, and a larger k indicates that \mathcal{X}_i^k is coarser. As shown in Figure 2, time series with multiple downsampling intervals (i.e., $k = 1, 2, 4$) is treated as the input data for training. To better utilize the small-loss criterion for clean label selection, each time series sample performs cross-scale fusion from fine to coarse (i.e., $A \rightarrow B, B \rightarrow C$) in the embedding space, which is mathematically defined as:

$$v_i^k = f\left(r_i^k \parallel v_i^{k-t} \parallel (r_i^k - v_i^{k-t}) \parallel (r_i^k \cdot v_i^{k-t})\right), \quad (2)$$

where r_i^k represents the single-scale embedding acquired by learning \mathcal{X}_i^k through an encoder. Meanwhile, v_i^k (or v_i^{k-t}) denotes the embedding of the time series X_i^k (or X_i^{k-t}) after performing cross-scale fusion. Here, t denotes the interval between adjacent downsampling ratios, and \parallel signifies the concatenation of two vectors to form a new vector. Notably, when $k = 1$, we employ the single-scale for classification training, resulting in $v_i^k = r_i^k$. By combining $(r_i^k - v_i^{k-t})$ and $(r_i^k \cdot v_i^{k-t})$ for vector concatenation, v_i^k can capture more nuanced discriminative information between r_i^k and v_i^{k-t} than that of simply concatenating r_i^k with v_i^{k-t} . The function $f(\cdot)$ represents a two-layer nonlinear network mapping function for fusing information of r_i^k and v_i^{k-t} . Additionally, v_i^k has the same dimension as r_i^k and serves as the input data for the multi-scale embedding graph learning process.

3.4 Multi-scale Embedding Graph Learning for Noisy Label Correction

We now present the multi-scale embedding graph learning module for correcting noisy labels. This module incorporates selected clean labels using label propagation theory. The process consists of two stages: graph construction and noisy label correction.

Graph Construction. It is assumed that the set of cross-fusion embeddings obtained from a batch of time series is defined as $V = \{v_1^k, v_2^k, \dots, v_M^k\}$, where M is the batch size. Intuitively, samples close to each other in the feature space have a high probability of belonging to the same class. However, in label-noise learning, v_i^k obtained from the current iterative training of the model may have unstable information, resulting in large deviations in the information of the nearest-neighbor samples of v_i^k . To address this issue, the proposed approach performs a momentum update [49] on v_i^k during training, which is defined as:

$$\bar{v}_i^k[e] = \alpha v_i^k[e] + (1 - \alpha) \bar{v}_i^k[e - 1], \quad (3)$$

where e is the current training epoch and α denotes the momentum update parameter.

The multi-scale embeddings nearest-neighbor graph can be created by using Euclidean distance among different \bar{v}_i^k . A common approach is the use of the Gaussian similarity function [45] to obtain the nearest-neighbor graph edge weight, which is defined as:

$$W_{ij} = \exp\left(-\frac{1}{2}d\left(\frac{\bar{v}_i^k}{\sigma}, \frac{\bar{v}_j^k}{\sigma}\right)\right), \quad (4)$$

where $d(\cdot)$ is the Euclidean distance function and σ is a fixed parameter. $W \in R^{M \times M}$ is a symmetric adjacency matrix, and the element W_{ij} denotes the nearest-neighbor edge weight between the embedding v_i^k and v_j^k (note that larger values indicate closer proximity). Then, W is normalized based on the graph laplacians [50] to obtain $Q = D^{-1/2}WD^{-1/2}$, where $D = \text{diag}(W1_n)$ is a diagonal matrix. Specifically, the K neighbors with the largest values in each row of Q are employed to create the nearest-neighbor graph. It is noteworthy that the embeddings in each mini-batch are utilized to generate the nearest-neighbor graph, thus obtaining Q within short computational time.

Noisy Label Correction. Specifically, small training loss samples acquired by DNNs in the early training period can be considered as clean samples, while samples with large training losses are considered as noisy ones [12, 14]. The above learning pattern of DNNs has been mathematically validated [11] (see Appendix A for details). Under this criterion, prior studies [12, 13, 23] have typically employed samples with small losses after a e_{warm} warm-up training as clean labels. Following [29], we extend the small-loss sample selection process to operate within each class, thereby enhancing the overall quality of the chosen clean labels. In our method, samples chosen with clean labels are considered labeled data, whereas unselected samples are treated as unlabeled data.

We utilize clean samples selected from time series at different scales in a cross-teaching manner (as in Figure 2). This could explore complementary information from different scale fusion embeddings to deal with noisy labels. It is supposed that there is a corresponding one-hot encoding matrix $Y \in R^{M \times C}$ ($Y_{ij} \in \{0, 1\}$) for the cross-fusion embeddings V . If y_i is identified as a clean label, we employ y_i to set Y_i as a one-hot encoded label. Otherwise, all the elements in Y_i are identified as zero. Through Y , the pseudo-label of each node in the nearest-neighbor graph Q can be obtained in an iterative way based on the label propagation theory. The specific solution formula is defined as:

$$F_{t+1} = \beta Q F_t + (1 - \beta) Y, \quad (5)$$

where $F_t \in R^{M \times C}$ denotes the predicted pseudo-label of the t -th iteration and $\beta \in (0, 1)$ is a hyperparameter. Naturally, F_t has a closed-form solution [24] defined as follows:

$$\mathcal{F} = (I - \beta Q)^{-1} Y, \quad (6)$$

where $\mathcal{F} \in R^{M \times C}$ is the final pseudo-labels and I denotes the identity matrix. Finally, the corrected label obtained for an unselected large-loss sample X_i is defined as:

$$y_i = \arg \max_c \mathcal{F}_i^c, \quad (7)$$

However, \mathcal{F} is the estimated pseudo-labels, which inevitably contain some incorrect labels. To address this issue, two strategies are used to improve the quality of pseudo-labels in \mathcal{F} . For the first strategy, the model continues training e_{update} epochs by using small-loss samples after e_{warm} epochs warm-up training to improve the robustness of the multi-scale embeddings. Then, the noisy label correction is performed after $(e_{warm} + e_{update})$ epoch. For the second strategy, a dynamic threshold $\varphi_e(c) = \frac{\delta_e(c)}{\max(\delta_e)} \gamma$ is utilized for each class [51] to select the pseudo-labels with a high confidence for noisy label correction, where $\delta_e(c)$ is the number of labeled samples contained in class c in the e -th epoch, and γ is a constant threshold.

Overall Training. Finally, each encoder utilizes the selected clean samples in combination with multi-scale embedding graph learning to perform noisy label correction for unselected large-loss samples. Combining the training data of the selected clean labels and those of corrected labels, the proposed Scale-teaching paradigm utilizes cross-entropy for time-series label-noise learning. Please refer to Algorithm 1 in the Appendix for the specific pseudo-code of Scale-teaching.

4 Experiments

4.1 Experiment Setup

Datasets. We use three time series benchmarks (four individual large datasets [3, 52, 53], UCR 128 archive [22], and UEA 30 archive [54]) for experiments. Among the four individual large datasets, HAR [52] and UniMiB-SHAR [3] are human activity recognition scenarios; FD-A [53] is the mechanical fault diagnosis scenario; Sleep-EDF [52] belongs to the sleep stage classification scenario. The UCR archive [22] contains 128 univariate time series datasets from different real-world scenarios. The UEA archive [54] contains 30 multivariate time series datasets from real-world scenarios. For details on the above datasets, please refer to Appendix B. Since all the datasets in three time series benchmarks are correctly labeled, we utilize a label transformation matrix T to add noises to the original correct labels [4], where T_{ij} is the probability of label i being flipped to j . We use three types of noisy labels for evaluations, namely Symmetric (Sym) noise, Asymmetric (Asym) noise, and Instance-dependent (Ins) noise. Symmetric (Asymmetric) noise randomly replaces a true label with other labels with an equal (unequal) probability. Instance noise [55] means that the noisy label is instance-dependent. Like [4, 12, 23], we use the test set with correct labels for evaluations.

Table 1: Test classification accuracy results compared with baselines on three time series benchmarks. The best results are **bold**, and the second best results are underlined. When P-value < 0.05, it indicates that the performance of Scale-teaching is statistically significant than the baseline.

Dataset	Noise Ratio	Metric	Standard	Mixup	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
Four individual large datasets	Sym 20%	Avg Rank	4.75	4.75	4.50	7.50	6.50	4.50	<u>2.50</u>	1.00
	Sym 50%	Avg Rank	4.75	4.50	4.75	7.25	5.75	4.50	<u>3.25</u>	1.25
	Asym 40%	Avg Rank	5.00	5.50	3.75	7.50	5.75	4.00	<u>3.25</u>	1.00
	Ins 40%	Avg Rank	4.75	4.25	4.25	7.25	6.00	4.75	<u>3.50</u>	1.00
UCR 128 archive	Sym 20%	Avg Rank	4.15	4.33	3.61	7.50	6.16	3.48	3.54	3.02
		P-value	1.90E-04	4.06E-05	1.90E-03	1.49E-34	1.70E-17	3.04E-03	8.57E-03	-
	Sym 50%	Avg Rank	4.31	4.57	4.05	6.43	5.89	<u>3.56</u>	3.86	3.11
		P-value	3.15E-05	1.70E-05	4.02E-04	7.48E-19	1.22E-15	1.40E-02	4.93E-03	-
	Asym 40%	Avg Rank	4.38	4.80	3.93	6.91	5.91	3.30	3.67	2.95
		P-value	1.62E-05	3.53E-07	6.10E-04	1.93E-23	9.82E-14	1.89E-02	2.24E-02	-
	Ins 40%	Avg Rank	4.05	4.52	4.02	7.04	6.18	3.30	3.77	2.95
		P-value	1.43E-05	1.81E-06	2.43E-04	9.81E-26	2.36E-17	3.27E-02	1.54E-02	-
UEA 30 archive	Sym 20%	Avg Rank	5.03	5.20	3.83	6.37	4.77	3.73	4.00	2.73
		P-value	6.61E-04	3.33E-04	2.69E-02	2.37E-05	1.14E-02	2.63E-02	3.93E-02	-
	Sym 50%	Avg Rank	5.17	5.73	4.23	6.23	3.93	3.83	4.30	2.43
		P-value	2.98E-04	7.40E-05	1.59E-02	9.35E-05	1.67E-02	1.08E-02	3.75E-02	-
	Asym 40%	Avg Rank	5.60	4.77	4.40	6.13	4.20	4.00	3.97	2.73
		P-value	3.81E-03	6.17E-03	1.63E-02	9.33E-05	1.36E-02	2.62E-02	3.88E-02	-
	Ins 40%	Avg Rank	5.20	4.77	4.33	6.60	4.27	4.20	3.77	2.60
		P-value	6.08E-04	2.92E-03	1.20E-02	2.55E-05	5.52E-03	1.08E-02	3.47E-02	-

Baselines. We select seven methods for comparative analyses, namely 1) Standard: direct training of the model using cross-entropy with all noisy labels; 2) Mixup [56]; 3) Co-teaching [12]; 4) FINE [35]; 5) SREA [4]; 6) SELC [37]; and 7) CULCU [23]. Among them, Standard, Mixup, and Co-teaching are the benchmark methods for label-noise learning. FINE, SELC, and CULCU are the state-of-the-art methods that do not need to focus on data types, and SREA is the state-of-the-art method in time series domain. In addition, for fair comparisons, all the baselines and the proposed Scale-teaching paradigm use the same encoder and classifier. We focus on the ability of different label-noise learning paradigms to cope with time series noise labels, rather than the classification performance achieved by using fully correct labels. Hence, considering the trade-off between the running time and classification performance, we choose FCN [57] as the encoder of Scale-teaching. For more details of baselines, please refer to Appendix C.

Implementation Details. Based on the experience [19, 44] in time series modeling, we utilize three different sampling intervals 1, 2, 4 as the input multi-scale series data for Scale-teaching. We use Adam as the optimizer. The learning rate is set to 1e-3, the maximum batch size is set to 256, and the maximum epoch is set to 200. e_{warm} is set to 30 and e_{update} is set to 90. α in Eq. 3 is set to 0.9, σ in Eq. 4 is set to 0.25, β in Eq. 5 is set to 0.99, the largest neighbor K is set to 10, and γ is set to 0.99. In addition, following the parameter settings suggested in [23], we linearly decay the learning rate to zero from the 80-th epoch to 200-th epoch. For a comprehensive understanding of the hyperparameter selection and the implementation of the small-loss criterion applied to Scale-teaching, please consult Appendix C. To reduce random errors, we utilize the mean test classification accuracy of the last five epochs of the model on the test set as experimental results. All the experiments are independently conducted five times with five different seeds, and the average classification accuracy and rank are reported. Finally, we build our model using PyTorch 1.10 platform with 2 NVIDIA GeForce RTX 3090 GPUs. Our implementation of Scale-teaching is available at <https://github.com/qianlima-lab/Scale-teaching>.

4.2 Main Results

We evaluate each time series benchmark using four noise ratios, Sym 20%, Sym 50%, Asym 40%, and Ins 40%. Due to space constraints, we only give the average ranking of all the methods on each benchmark in Table 1. Please refer to Appendix D for the specific test classification accuracies. Besides, for UCR 128 and UEA 30 archives, we use the Wilcoxon signed rank test (P-value) [58] to analyze the classification performance of baselines. As shown in Table 1, the proposed Scale-teaching paradigm achieves the best Avg Rank in all the cases. It is found that Mixup [56] and FINE [35] perform worse than the Standard method in most cases. For Mixup, the complex dynamic properties

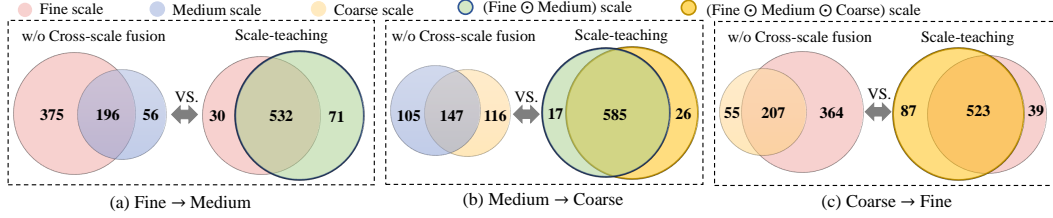


Figure 3: Venn diagram of the average number of correctly classified samples for the different scale sequences of UCR 128 archive with Sym 20% noisy labels. The numbers in the figure indicate the complements and intersections of classification results at different scales.

Table 2: The test classification accuracy (%) results of different scale classifiers on UCR 128 archive. The best results are **bold**, and the second best results are underlined. When P-value < 0.05, it indicates that the performance of Scale-teaching’s coarse scale classifier is significant than other classifiers.

Method		w/o Cross-scale fusion			Scale-teaching		
Noise Ratio	Metric	Fine	Medium	Coarse	Fine	Medium	Coarse
Sym 20%	Avg Acc	65.13	30.11	28.17	59.67	<u>68.17</u>	68.70
	Avg Rank	2.38	5.09	5.37	3.20	<u>2.17</u>	2.11
	P-value	1.89E-03	2.85E-37	2.07E-40	1.58E-09	3.74E-02	-
Asym 40%	Avg Acc	49.61	29.01	28.87	47.75	<u>51.93</u>	52.87
	Avg Rank	2.64	4.78	4.75	3.01	<u>2.45</u>	2.27
	P-value	1.94E-03	6.78E-25	1.59E-27	1.80E-07	2.80E-02	-

of the original time series are destroyed probably due to the mixture of two different time series mechanisms. FINE uses embeddings of the input data to select clean labels. Although FINE achieves advanced classification performance for image data, it is difficult to be used directly for time series data because its discriminative patterns are easily distorted by external noises. SREA [4] has a good performance on the UEA 30 archive, while it performs poorly on the other benchmarks. Meanwhile, Co-teaching [12], SELC [37], and CULCU [23] are more robust against time series noisy labels in different cases, further indicating that the small-loss criterion is also applicable to time series.

4.3 Multi-scale Analysis

To explain the multi-scale mechanism in the Scale-teaching paradigm, we add an ablation study based on Scale-teaching (w/o cross-scale fusion). We select the UCR 128 archive to analyze the classification results obtained by the fine, medium, and coarse scale classifiers. As shown in Figure 3, the classification results of different scale sequences have evident complementary information. Scale-teaching can effectively use complementary information between cross-scale to obtain more robust embeddings and clean labels. In response to the tendency of the coarse scale to ignore discriminative patterns in fine scale (please see Table 2), our proposed cross-scale fusion mechanism can effectively improve the classification performance of medium and coarse scales while retaining complementarity. Please refer to Appendix E for the specific classification results of Figure 3 and Table 2. In Appendix E, we also analyze the order and size of the downsampled input scale sequence for Scale-teaching.

Scale-teaching utilizes multi-scale embeddings to generate the nearest-neighbor graph, and uses clean labels selected for noisy label correction. To explore the distribution of different classes of embeddings, we employ t-SNE [59] for dimensionality reduction visualization. Specifically, we utilize the UniMiB-SHAR dataset containing Sym 20% noisy labels for visualization. As shown in Figure 4, we find that the embeddings learned by Scale-teaching are more discriminative across classes than the Standard and CULCU methods that use a single scale series for training. The above results suggest that Scale-teaching can effectively exploit the complementary information between different scales, prompting the learned embeddings to be more discriminative between classes. In addition, we choose the FD-A dataset for t-SNE visualization, and please refer to Appendix E.

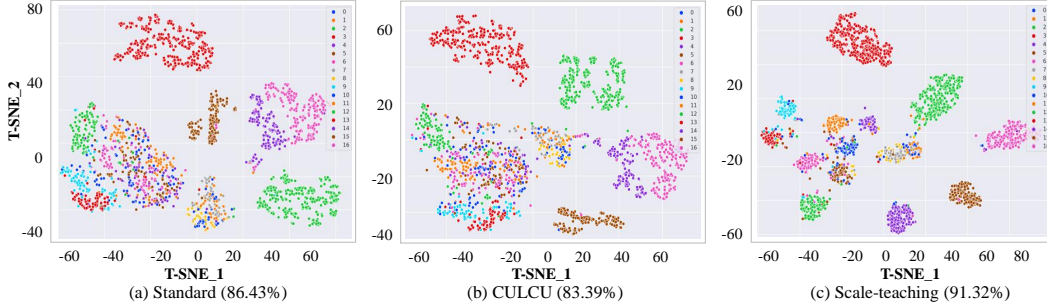


Figure 4: t-SNE visualization of the learned embeddings on the UnimiB-SHAR dataset with Sym 20% noisy labels (values in parentheses are the test classification accuracies).

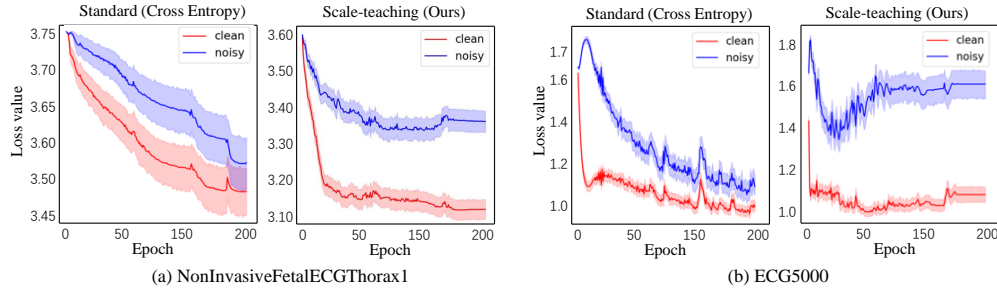


Figure 5: The change of loss values for clean and noisy time series samples under Aysm 40% noise labels. The solid line and shading indicate the mean and standard deviation loss values of all clean (or noisy) training samples within each epoch.

4.4 Small-loss Analysis

To analyze the application of the small-loss criterion to time series data, we visualize the change of loss values for ground-truth clean and noisy time series samples during training. Specifically, Figure 5 shows the change in loss values of the models trained by the Standard method and Scale-teaching on two UCR time series datasets. When the model is trained with the Standard method, differences can be found in the loss values of clean and noisy samples in the network early training, especially in Figure 5 (b). The Standard method makes the model gradually fit the noisy samples as the training proceeds, while Scale-teaching improves the ability of the model to handle noisy labels. To further prove its effectiveness, we selected two other UCR datasets for the loss value change analysis, which have the same pattern as Figure 5. Also, we report the HAR and UniMiB-SHAR dataset’s loss value probability distributions of clean and noisy samples. For more details, please refer to Appendix F.

4.5 Ablation Study

To verify the robustness of each module in Scale-teaching, the ablation experiments have been conducted in the HAR and UniMiB-SHAR datasets, and the results are shown in Table 3. Specifically, (1) **w/o cross-scale fusion**: the cross-scale embedding fusion from fine to coarse mechanism is ablated; (2) **only single scale**: only the original time series is used for training; (3) **w/o graph learning**: the multi-scale embedding graph learning module for noisy label correction is ablated; (4) **w/o moment**: the embedding momentum update mechanism (Eq. 3) is ablated; (5) **w/o dynamic threshold**: using a dynamic threshold to select high-quality propagation pseudo-labels is ablated.

As shown in Table 3, the cross-scale fusion strategy (w/o cross-scale fusion) and the clean labels cross-teaching mechanism (only single scale) can effectively improve the classification performance of Scale-teaching, especially on the UniMiB-SHAR dataset with a large number of classes. Meanwhile, in terms of label correction based on multi-scale embedding graph learning, the results of the corresponding ablation module show that improving the stability of embedding (w/o moment) and

Table 3: The test classification accuracy (%) results of ablation study (values in parentheses denote drop accuracy).

Method	HAR		UniMiB-SHAR	
	Sym 50%	Asym 40%	Sym 50%	Asym 40%
Scale-teaching	90.17	89.62	81.31	70.68
w/o cross-scale fusion	88.47 (-1.70)	87.64 (-1.98)	73.32 (-7.99)	61.62 (-9.06)
only single scale	89.01 (-1.06)	88.11 (-1.51)	69.89 (-11.42)	60.32 (-10.36)
w/o graph learning	88.06 (-2.11)	87.65 (-1.97)	79.72 (-1.59)	68.87 (-1.81)
w/o moment	89.76 (-0.41)	88.76 (-0.86)	80.57 (-0.74)	69.85 (-0.83)
w/o dynamic threshold	89.12 (-1.05)	88.75 (-0.87)	77.42 (-3.89)	69.53 (-1.15)

Table 4: The test classification accuracy (%) results on four individual large datasets without noisy labels. The best results are **bold**, and the second best results are underlined.

Dataset	Standard	Mixup	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
HAR	93.29	95.42	93.77	93.13	93.02	93.76	94.75	<u>94.72</u>
UniMiB-SHAR	89.14	84.84	88.24	88.14	65.51	89.28	<u>89.46</u>	93.61
FD-A	99.93	99.91	99.96	68.22	90.25	99.82	99.95	99.96
Sleep-EDF	84.93	84.67	<u>85.37</u>	84.62	79.42	84.82	85.54	85.34

selecting high-quality pseudo-labels (w/o dynamic threshold) can effectively improve the performance of label correction based on graph learning.

Furthermore, we select the four individual large datasets without noisy labels for evaluation. As shown in Table 4, Scale-teaching’s classification performance is still better than most baselines. It’s worth mentioning that SREA [4] employs an unsupervised time series reconstruction loss as an auxiliary task, which reduces the model’s classification performance without noisy labels. We also provide the corresponding test classification results for Tables 3 and 4 under the F1-score metric in Appendix G. Additionally, we find the running time of Scale-teaching, which is faster than FINE, SREA and CULCU for datasets with a larger number of samples or longer length of the sequence. We further analyze the classification performance of the proposed Scale-teaching paradigm and time series classification methods [15, 60] in Appendix G.

5 Conclusions

Limitations. The input scales of our proposed Scale-teaching paradigm can only select a fixed number of scales for training, and the running time will increase as the number of scales increases.

Conclusion. In this paper, we propose a deep learning paradigm for time-series classification with noisy labels called Scale-teaching. Experiments on the three time series benchmarks show that the Scale-teaching paradigm can utilize the multi-scale properties of time series to effectively handle noisy labels. Comprehensive analyses on multi-scale and ablation studies demonstrate the robustness of the Scale-teaching paradigm. In the future, we will explore the design of scale-adaptive time-series label-noise learning models.

Acknowledgments

We thank the anonymous reviewers for their helpful feedbacks. We thank Professor Eamonn Keogh and all the people who have contributed to the UCR 128 archive, UEA 30 archive, and the four large individual time series classification datasets. The work described in this paper was partially funded by the National Natural Science Foundation of China (Grant Nos. 62272173, 61872148, 62206041), the Natural Science Foundation of Guangdong Province (Grant Nos. 2022A1515010179, 2019A1515010768), the Science and Technology Planning Project of Guangdong Province (Grant No. 2023A0505050106), the Fundamental Research Funds for the Central Universities under grants DUT22RC(3)015. The authors would like to thank Siying Zhu, Huawen Feng, Yu Chen, and Junlong Liu from SCUT for their review and helpful suggestions.

References

- [1] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data mining and knowledge discovery*, 33(4):917–963, 2019.
- [2] Mingyue Cheng, Qi Liu, Zhiding Liu, Zhi Li, Yucong Luo, and Enhong Chen. Formertime: Hierarchical multi-scale representations for multivariate time series classification. *arXiv preprint arXiv:2302.09818*, 2023.
- [3] Gentry Atkinson and Vangelis Metsis. Tsar: a time series assisted relabeling tool for reducing label noise. In *The 14th PErvasive Technologies Related to Assistive Environments Conference*, pages 203–209, 2021.
- [4] Andrea Castellani, Sebastian Schmitt, and Barbara Hammer. Estimating the electrical power output of industrial devices with end-to-end time-series classification in the presence of label noise. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 469–484. Springer, 2021.
- [5] Benoît Frénay, Gaël de Lannoy, and Michel Verleysen. Label noise-tolerant hidden markov models for segmentation: application to eegs. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 455–470. Springer, 2011.
- [6] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [7] Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. Robust training under label noise by over-parameterization. *arXiv preprint arXiv:2202.14026*, 2022.
- [8] Gentry Atkinson and Vangelis Metsis. A survey of methods for detection and correction of noisy labels in time series data. In *IFIP International Conference on Artificial Intelligence Applications and Innovations*, pages 479–493. Springer, 2021.
- [9] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *International conference on machine learning*, pages 2304–2313. PMLR, 2018.
- [10] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [11] Xian-Jin Gui, Wei Wang, and Zhang-Hao Tian. Towards understanding deep learning from noisy labels with small-loss criterion. *arXiv preprint arXiv:2106.09291*, 2021.
- [12] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31, 2018.
- [13] Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In *International Conference on Machine Learning*, pages 7164–7173. PMLR, 2019.
- [14] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020.
- [15] Patrick Schäfer. The boss is concerned with time series classification in the presence of noise. *Data Mining and Knowledge Discovery*, 29:1505–1530, 2015.
- [16] Zhicheng Cui, Wenlin Chen, and Yixin Chen. Multi-scale convolutional neural networks for time series classification. *arXiv preprint arXiv:1603.06995*, 2016.
- [17] Patrick Schäfer and Ulf Leser. Fast and accurate time series classification with weasel. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 637–646, 2017.

- [18] Shiyu Chang, Yang Zhang, Wei Han, Mo Yu, Xiaoxiao Guo, Wei Tan, Xiaodong Cui, Michael Witbrock, Mark A Hasegawa-Johnson, and Thomas S Huang. Dilated recurrent neural networks. *Advances in neural information processing systems*, 30, 2017.
- [19] Zipeng Chen, Qianli Ma, and Zhenxi Lin. Time-aware multi-scale rnns for time series modeling. In *IJCAI*, pages 2285–2291, 2021.
- [20] Wensi Tang, Guodong Long, Lu Liu, Tianyi Zhou, Michael Blumenstein, and Jing Jiang. Omni-scale cnns: a simple and effective kernel size configuration for time series classification. In *International Conference on Learning Representations*, 2022.
- [21] Víctor Campos, Brendan Jou, Xavier Giró-i Nieto, Jordi Torres, and Shih-Fu Chang. Skip rnn: Learning to skip state updates in recurrent neural networks. *arXiv preprint arXiv:1708.06834*, 2017.
- [22] Hoang Anh Dau, Anthony Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn Keogh. The ucr time series archive. *IEEE/CAA Journal of Automatica Sinica*, 6(6):1293–1305, 2019.
- [23] Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama. Sample selection with uncertainty of losses for learning with noisy labels. In *International Conference on Learning Representations*, 2022.
- [24] Dengyong Zhou, Olivier Bousquet, Thomas Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. *Advances in neural information processing systems*, 16, 2003.
- [25] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 322–330, 2019.
- [26] Xiong Zhou, Xianming Liu, Junjun Jiang, Xin Gao, and Xiangyang Ji. Asymmetric loss functions for learning with noisy labels. In *International conference on machine learning*, pages 12846–12856. PMLR, 2021.
- [27] De Cheng, Yixiong Ning, Nannan Wang, Xinbo Gao, Heng Yang, Yuxuan Du, Bo Han, and Tongliang Liu. Class-dependent label-noise learning with cycle-consistency regularization. *Advances in Neural Information Processing Systems*, 35:11104–11116, 2022.
- [28] Kilian Fatras, Bharath Bhushan Damodaran, Sylvain Lobry, Remi Flamary, Devis Tuia, and Nicolas Courty. Wasserstein adversarial regularization for learning with label noise. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [29] Nazmul Karim, Mamshad Nayeem Rizve, Nazanin Rahnavard, Ajmal Mian, and Mubarak Shah. Unicon: Combating label noise through uniform selection and contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9676–9686, 2022.
- [30] Deep Patel and PS Sastry. Adaptive sample selection for robust learning under label noise. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3932–3942, 2023.
- [31] Bo Han, Jiangchao Yao, Gang Niu, Mingyuan Zhou, Ivor Tsang, Ya Zhang, and Masashi Sugiyama. Masking: A new perspective of noisy supervision. *Advances in neural information processing systems*, 31, 2018.
- [32] Zhaowei Zhu, Jialu Wang, and Yang Liu. Beyond images: Label noise transition matrix estimation for tasks with lower-quality features. In *International Conference on Machine Learning*, pages 27633–27653. PMLR, 2022.
- [33] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. *Advances in neural information processing systems*, 33:20331–20342, 2020.

- [34] Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In *International conference on artificial intelligence and statistics*, pages 4313–4324. PMLR, 2020.
- [35] Taehyeon Kim, Jongwoo Ko, JinHwan Choi, Se-Young Yun, et al. Fine samples for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34:24137–24149, 2021.
- [36] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *International conference on machine learning*, pages 312–321. PMLR, 2019.
- [37] Yangdi Lu and Wenbo He. Selc: Self-ensemble label correction improves learning with noisy labels. *arXiv preprint arXiv:2205.01156*, 2022.
- [38] Lifeng Shen, Zhuocong Li, and James Kwok. Timeseries anomaly detection using temporal hierarchical one-class network. *Advances in Neural Information Processing Systems*, 33:13016–13026, 2020.
- [39] Qianli Ma, Zhen Liu, Zhenjing Zheng, Ziyang Huang, Siying Zhu, Zhongzhong Yu, and James T Kwok. A survey on time-series pre-trained models. *arXiv preprint arXiv:2305.10716*, 2023.
- [40] Zhiwei Wang, Zhengzhang Chen, Jingchao Ni, Hui Liu, Haifeng Chen, and Jiliang Tang. Multi-scale one-class recurrent neural networks for discrete event sequence anomaly detection. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 3726–3734, 2021.
- [41] Junchen Ye, Zihan Liu, Bowen Du, Leilei Sun, Weimiao Li, Yanjie Fu, and Hui Xiong. Learning the evolutionary and multi-scale graph structure for multivariate time series forecasting. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2296–2306, 2022.
- [42] Amin Shabani, Amir Abdi, Lili Meng, and Tristan Sylvain. Scaleformer: iterative multi-scale refining transformers for time series forecasting. *arXiv preprint arXiv:2206.04038*, 2022.
- [43] Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. Micn: Multi-scale local and global context modeling for long-term series forecasting. In *The Eleventh International Conference on Learning Representations*, 2023.
- [44] Tian Zhou, Ziqing Ma, Qingsong Wen, Liang Sun, Tao Yao, Rong Jin, et al. Film: Frequency improved legendre memory model for long-term time series forecasting. *arXiv preprint arXiv:2205.08897*, 2022.
- [45] Fei Wang and Changshui Zhang. Label propagation through linear neighborhoods. In *Proceedings of the 23rd international conference on Machine learning*, pages 985–992, 2006.
- [46] Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. *arXiv preprint arXiv:1805.10002*, 2018.
- [47] Konstantinos Kamnitsas, Daniel Castro, Loic Le Folgoc, Ian Walker, Ryutaro Tanno, Daniel Rueckert, Ben Glocker, Antonio Criminisi, and Aditya Nori. Semi-supervised learning via compact latent space clustering. In *International conference on machine learning*, pages 2459–2468. PMLR, 2018.
- [48] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5070–5079, 2019.
- [49] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016.
- [50] Fan RK Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997.

- [51] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34:18408–18419, 2021.
- [52] Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee Keong Kwoh, Xiaoli Li, and Cuntai Guan. Time-series representation learning via temporal and contextual contrasting. *arXiv preprint arXiv:2106.14112*, 2021.
- [53] Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. Self-supervised contrastive pre-training for time series via time-frequency consistency. *arXiv preprint arXiv:2206.08496*, 2022.
- [54] Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018. *arXiv preprint arXiv:1811.00075*, 2018.
- [55] Jiacheng Cheng, Tongliang Liu, Kotagiri Ramamohanarao, and Dacheng Tao. Learning with bounded instance and label-dependent label noise. In *International Conference on Machine Learning*, pages 1789–1799. PMLR, 2020.
- [56] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [57] Zhiguang Wang, Weizhong Yan, and Tim Oates. Time series classification from scratch with deep neural networks: A strong baseline. In *2017 International joint conference on neural networks (IJCNN)*, pages 1578–1585. IEEE, 2017.
- [58] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7:1–30, 2006.
- [59] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [60] Angus Dempster, François Petitjean, and Geoffrey I Webb. Rocket: exceptionally fast and accurate time series classification using random convolutional kernels. *Data Mining and Knowledge Discovery*, 34(5):1454–1495, 2020.
- [61] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra Perez, and Jorge Luis Reyes Ortiz. A public domain dataset for human activity recognition using smartphones. In *Proceedings of the 21th international European symposium on artificial neural networks, computational intelligence and machine learning*, pages 437–442, 2013.
- [62] Daniela Micucci, Marco Mobilio, and Paolo Napoletano. Unimib shar: A dataset for human activity recognition using acceleration data from smartphones. *Applied Sciences*, 7(10):1101, 2017.
- [63] Christian Lessmeier, James Kuria Kimotho, Detmar Zimmer, and Walter Sextro. Condition monitoring of bearing damage in electromechanical drive systems by using motor current signals of electric motors: A benchmark data set for data-driven classification. In *PHM Society European Conference*, volume 3, 2016.
- [64] Ralph G Andrzejak, Klaus Lehnertz, Florian Mormann, Christoph Rieke, Peter David, and Christian E Elger. Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: Dependence on recording region and brain state. *Physical Review E*, 64(6):061907, 2001.
- [65] Hassan Ismail Fawaz, Benjamin Lucas, Germain Forestier, Charlotte Pelletier, Daniel F Schmidt, Jonathan Weber, Geoffrey I Webb, Lhassane Idoumghar, Pierre-Alain Muller, and François Petitjean. Inceptiontime: Finding alexnet for time series classification. *Data Mining and Knowledge Discovery*, 34(6):1936–1962, 2020.

Supplementary Material: Scale-teaching: Robust Multi-scale Training for Time Series Classification with Noisy Labels

A Small-loss Criterion

DNNs have been widely known to first learn simple and generalized patterns, which is achieved by learning clean data. After that, the networks gradually overfit noisy ones. In other words, when we train a model with a dataset containing incorrectly labeled samples, we can consider the samples with small training losses as clean ones and use them to update the model. Formally, let f^* be the target concept which determines the true label of x and model $g^* = g(x; \Theta^*)$ minimizing the expected loss, i.e.,

$$\Theta^* = \arg \min_{\Theta} \mathbb{E}_{(x, \tilde{y})} [\ell_{CE}(g(x; \Theta), \tilde{y})]. \quad (8)$$

Then, the small-loss criterion can be stated as follows[1]:

Theorem 1. *Suppose g is ϵ -close to g^* , i.e., $\|g - g^*\|_{\infty} = \epsilon$, for two examples (x_1, \tilde{y}) and (x_2, \tilde{y}) , assume $f^*(x_1) = \tilde{y}$ and $f^*(x_2) \neq \tilde{y}$, if T satisfies the diagonally-dominant condition $T_{ii} > \max\{\max_{j \neq i} T_{ij}, \max_{j \neq i} T_{ji}\}, \forall i$, and $\epsilon < \frac{1}{2} \cdot (T_{\tilde{y}\tilde{y}} - T_{f^*(x_2)\tilde{y}})$, then $\ell_{CE}(g(x_1), \tilde{y}) < \ell_{CE}(g(x_2), \tilde{y})$.*

The work [11] provides the proof of this theorem. It shows that during training, the model can select clean samples according to the loss values. The reason is that the loss values of clean samples among the samples with the same observed labels are smaller. It is worth noting that the theorem is under the assumption of the class-dependent noise type and requires the transition matrix to satisfy the diagonally-dominant condition. Additionally, the finite data may also make the conditions of the theorem difficult to hold because the model g may be far away from g^* .

B Dataset Information

To evaluate the robustness of our proposed Scale-teaching and baselines on the time-series label-noise learning task, we selected three benchmark time-series datasets for experimental analysis.

B.1 Four individual large datasets

The statistical information of the four individual time series datasets is shown in Table 5. And the specific dataset information is as follows:

Human Activity Recognition (HAR)

The HAR dataset [52, 61] is collected from 30 students performing six human actions (i.e., walking, walking upstairs, downstairs, standing, sitting, and lying down) by wearing sensors.

University of Milano Bicocca Smartphone-based Human Activity Recognition (UniMiB SHAR)

The UniMiB SHAR dataset [3, 62] is human activity information collected at a sampling rate of 50 Hz from volunteers with a smartphone with an accelerometer sensor in the front pocket of their pants. Specifically, each accelerometer entry is labeled by specifying the type of ADL (e.g., walking, sitting, or standing) or the type of fall (e.g., forward, fainting, or backward).

Faulty Detection Condition A (FD-A)

The FD-A dataset [52, 63] is generated by an electromechanical drive system that monitors the condition of rolling bearings and detects their failure. Each rolling bearing can be classified into three categories: undamaged, inner damaged, and externally damaged.

Sleep Stage EEG Signal Classification (Sleep-EDF)

The Sleep-EDF dataset [52, 64] includes the whole night PSG sleep recordings, which contain five EEG sleep signal recordings: Wake (W), Non-rapid eye movement (N1, N2, N3), and Rapid Eye Movement (REM).

Table 5: A summary of four individual large time series datasets used in the experiments.

Dataset	# Train	# Test	Length	# Variables	# Classes
HAR	7352	2947	128	9	6
Sleep-EDF	25612	8910	3000	1	5
FD-A	8184	2728	5120	1	3
UniMiB-SHAR	9416	2354	453	1	17

B.2 UCR 128 Archive

The UCR time series archive [22] contains 128 univariate datasets and is widely used for classification in the time series mining community. Each UCR dataset includes a single training set and a single test set, and each time series sample has been z-normalized. In addition, we uniformly use the mean-imputation method to preprocess the datasets that contain missing values. For detailed information about UCR datasets, please refer to https://www.cs.ucr.edu/~eamonn/time_series_data_2018/.

B.3 UEA 30 Archive

The UEA time series archive [54] contains 30 multivariate datasets, mainly derived from Human Activity Recognition, Motion classification, ECG classification, EEG/MEG classification, Audio Spectra Classification, and other realistic scenarios. Each dataset contains a partitioned training set and a test set. In addition, we use the mean-imputation method to deal with datasets with missing values. For detailed information about UEA datasets, please refer to <https://www.timeseriesclassification.com/dataset.php>.

C Baselines

To analyze the performance and effectiveness of Scale-teaching on time-series label-noise learning, we selected seven baselines for comparative analysis. The specific information is as follows.

- Standard directly employs all samples in the training set containing noisy labels and performs supervised classification training using cross-entropy loss. Then, the trained model is used to make predictions on the test set.
- Mixup [56] trains a neural network on convex combinations of pairs of time series samples and their labels (whatever is clean or noisy). For the specific open source code, please refer to <https://github.com/facebookresearch/mixup-cifar10>.
- Co-teaching [12] trains two deep neural networks simultaneously, and lets them teach each other given every mini-batch with selected clean labels based on a small-loss criterion. For the specific open source code, please refer to <https://github.com/bhanML/Co-teaching>.
- FINE [35] utilizes a novel detector for clean label selection. Especially, FINE focus on each data point’s latent representation dynamics and measures the alignment between the latent distribution and each representation using the eigen decomposition of the data gram matrix. For the specific open source code, please refer to https://github.com/Kthyeon/FINE_official.
- SREA [4] employs a novel multi-task deep learning approach for time series noisy label correction that jointly trains a classifier and an autoencoder with a shared embedding representation. For the specific open source code, please refer to <https://github.com/Castel144/SREA>.
- SELC [37] utilizes a simple and effective method self-ensemble label correction (SELC) to progressively correct noisy labels and refine the model. For the specific open source code, please refer to <https://github.com/MacLLL/SELC>.
- CULCU [23] incorporates the uncertainty of losses by adopting interval estimation instead of point estimation of losses to select clean labels based on Co-teaching. CULCU has two

versions: CNLCU-S and CNLCU-H, where CNLCU-S uses soft labels for training and CNLCU-H uses hard labels for training. According to the original paper’s [23] experimental results, CNLCU-S has a better performance. Hence, we use CNLCU-S as a baseline. For the specific open source code, please refer to <https://github.com/xiaoboxia/CNLCU>.

Finally, based on the source code of the above baselines, we provide the reproduction source code of all baselines, as well as the source code of our proposed Scale-teaching (refer to Algorithm 1). For the specific open-source code, please refer to our GitHub repository <https://github.com/qianlima-lab/Scale-teaching>.

Our experiment contains 162 datasets. It would be time-consuming to perform hyperparameter selection for each dataset. Therefore, the hyperparameters of Scale-teaching are not carefully tuned for each dataset, and most of the hyperparameters are set based on the default hyperparameters of related works. The learning rate and maximum epoch are set based on the parameters of existing noise-label learning methods, such as FINE and CULCU. α in Eq. 3, σ in Eq. 4 and β in Eq. 5 are set based on the default hyperparameters of related label propagation works. e_{warm} is based on FINE settings. e_{update} , γ and batch size are based on manual empirical settings without specific hyperparameter analysis. The largest neighbor K is set based on human experience, and we had a simple test on several datasets, and found that a larger value of does not improve the classification performance, but instead increases the running time of the model.

For the implementation of small-loss criterion in Scale-teaching, we select small-loss samples within each class from the mini-batch data as clean labeled data. For stduies [12, 13], they use warm-up training to decrease $\lambda(e)$ from 1 to $1 - \eta$. $\lambda(e)$ denotes the selection ratio of small-loss samples within the mini-batch data without considering the difference of class, and η is the ratio of noise labels in the training set. Based on the above criterion, the current work [29] uses the Jensen-Shannon divergence to calculate difference d between the classification result p_i of sample \mathcal{X}_i^c and the observation label \hat{y}_i . Following [29], for each class c , we consider the observed label of \mathcal{X}_i^c as a clean label when the d of the training sample \mathcal{X}_i^c is less than d_{avg}^e after a e_{warm} warm-up training. d_{avg}^e denotes the average of ds of all the training samples when the epoch is e . We observed that using the Jensen-Shannon divergence method [29] and directly employing stduies [12, 13] for clean sample selection within each class have distinct strengths and weaknesses when applied to various time series datasets. In our study, we implemented the strategy of stduies [12, 13] for clean sample selection within each class on four individual large datasets and the UCR 128 archive. Meanwhile, the Jensen-Shannon divergence method [29] was applied to the UEA 30 archive for clean sample selection within each class.

D Details of Main Results

For the four individual large time series datasets, the specific classification results of our proposed Scale-teaching paradigm and baselines are shown in Table 6. For the UCR 128 archive, the specific classification results for all methods with different noise ratios are shown in Table 11 (Sym 20%), 12 (Sym 50%), 13 (Asym 40%), and 14 (Ins 40%). For the UEA 30 archive, the specific classification results for all methods at different noise ratios are shown in Tables 15 (Sym 20%), 16 (Sym 50%), 17 (Asym 40%), and 18 (Ins 40%). For layout and reading convenience, we only give the average classification accuracy for multiple runs of all methods without standard deviation on the UCR 128 archive and UEA 30 archive.

E Details of Multi-scale Results

To analyze the multi-scale mechanism in the Scale-teaching paradigm, we provide the classification performance of classifiers corresponding to fine, medium and coarse scales, as shown in Tables 19 and 21. And the classification results by ablation cross-scale fusion mechanism based on the Scale-teaching are shown in Tables 20 and 22. For the abbreviations in Tables 19, 20, 21 and 22, such as $a_t_b_f$, $b_f_c_t$, and $c_t_a_f$, where a denotes fine classifier, b denotes medium classifier, and c denotes coarse classifier, and t and f represent correct and incorrect classification results, respectively. For example, $a_t_b_f$ indicates the number of samples correctly predicted by the fine classifier and incorrectly predicted by the medium classifier. In addition, we provide t-SNE [59] visualization on the FD-A dataset with Sym 50% noisy labels (as in Figure 6) to explore the distribution of different classes of embeddings. Figure 6 shows that the cross-scale fusion mechanism in Scale-teaching for

Algorithm 1 The proposed Scale-teaching paradigm.

Input: encoders $[w_A, w_B, w_C]$, classifiers $[c_A, c_B, c_C]$, fine-scale series x_A , medium-scale series x_B , and coarse-scale series x_C

Output: $[w_A, w_B, w_C]$ and $[c_A, c_B, c_C]$

Note: For clarity, our analysis utilizes three distinct scales for training, but this approach can be extended to incorporate multiple scales.

- 1: **Step one:** Obtain single-scale embeddings r_A, r_B, r_C ;
 $r_A = w_A(x_A)$;
 $r_B = w_B(x_B)$;
 $r_C = w_C(x_C)$;
 - 2: **Step two:** Obtain cross-scale embeddings v_A, v_B, v_C ;
 $v_A = r_A$;
 $v_B = \text{Eq. 2}(r_B, v_A)$;
 $v_C = \text{Eq. 2}(r_C, v_B)$;
 - 3: **Step three:** Obtain clean labels y_A, y_B, y_C for cross-teaching training;
 $y_A = c_C(v_C)$ via small loss criterion;
 $y_B = c_A(v_A)$ via small loss criterion;
 $y_C = c_B(v_B)$ via small loss criterion;
 - 4: **Step four:** Obtain corrected labels yc_A, yc_B, yc_C for classification training;
 $yc_A = \text{Eq. 6}(v_A, y_A)$ via label propagation;
 $yc_B = \text{Eq. 6}(v_B, y_B)$ via label propagation;
 $yc_C = \text{Eq. 6}(v_C, y_C)$ via label propagation;
 - 5: **Step five:** Overall training;
Update encoder w_A and classifier c_A via cross-entropy loss(v_A, y_A & yc_A);
Update encoder w_B and classifier c_B via cross-entropy loss(v_B, y_B & yc_B);
Update encoder w_C and classifier c_C via cross-entropy loss(v_C, y_C & yc_C).
-

Table 6: The detailed test classification accuracy (%) compared with baselines on four individual large datasets (values in parentheses are standard deviations). The best results are in **bold**.

Dataset	Noise	Standard	Mixup	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
HAR	Sym 20%	92.13 (0.64)	92.52 (1.05)	92.28 (0.67)	92.15 (0.55)	92.53 (1.41)	92.88 (0.82)	92.66 (0.37)	93.93 (0.66)
	Sym 50%	83.99 (2.89)	76.75 (1.88)	89.90 (1.63)	88.42 (3.83)	91.38 (0.59)	90.37 (0.73)	89.91 (2.19)	90.17 (0.67)
	Asym 40%	75.59 (5.39)	66.91 (2.61)	87.67 (2.52)	83.87 (5.98)	88.98 (0.57)	87.67 (2.39)	87.22 (1.22)	89.62 (0.73)
	Ins 40%	83.56 (2.82)	73.86 (0.89)	90.98 (0.96)	90.77 (0.33)	91.25 (1.11)	91.02 (1.53)	91.15 (1.43)	91.58 (1.47)
UniMiB-SHAR	Sym 20%	87.07 (0.95)	82.13 (1.08)	80.54 (2.16)	26.63 (3.07)	51.48 (3.65)	68.52 (2.86)	82.80 (1.87)	90.69 (1.02)
	Sym 50%	79.37 (0.41)	77.77 (1.59)	66.33 (2.85)	18.92 (4.61)	47.62 (3.33)	67.65 (3.31)	66.36 (3.91)	81.31 (0.67)
	Asym 40%	63.59 (4.13)	66.32 (1.93)	60.25 (1.45)	19.18 (4.37)	51.16 (3.01)	55.65 (1.59)	60.45 (1.65)	70.68 (2.15)
	Ins 40%	55.83 (8.14)	56.97 (6.48)	54.09 (3.79)	11.18 (4.75)	51.5 (1.98)	54.62 (6.63)	53.90 (4.75)	71.14 (3.99)
FD-A	Sym 20%	98.89 (0.05)	99.78 (0.06)	99.83 (0.08)	78.13 (21.47)	89.92 (0.68)	99.67 (0.09)	99.85 (0.08)	99.93 (0.04)
	Sym 50%	96.63 (1.16)	98.73 (0.62)	99.04 (0.32)	70.65 (17.53)	82.18 (0.01)	98.59 (0.25)	99.06 (0.29)	99.38 (0.53)
	Asym 40%	96.12 (1.65)	93.50 (1.85)	97.06 (4.05)	61.04 (14.24)	90.23 (0.02)	98.24 (0.58)	98.91 (0.42)	99.55 (0.36)
	Ins 40%	99.36 (0.47)	99.55 (0.10)	99.51 (0.19)	67.81 (12.95)	88.63 (0.02)	99.36 (0.23)	99.53 (0.22)	99.82 (0.06)
Sleep-EDF	Sym 20%	85.01 (0.09)	84.31 (0.36)	84.81 (0.14)	81.21 (0.28)	72.79 (0.99)	84.32 (0.33)	85.23 (0.14)	85.56 (0.35)
	Sym 50%	83.58 (0.74)	83.61 (0.39)	83.39 (0.25)	78.17 (4.42)	72.78 (1.30)	83.06 (0.29)	84.02 (0.53)	84.59 (0.97)
	Asym 40%	79.62 (2.39)	77.40 (1.92)	82.87 (0.40)	64.77 (2.10)	72.23 (0.89)	82.50 (1.07)	83.05 (0.64)	83.87 (0.38)
	Ins 40%	84.35 (0.38)	84.25 (0.31)	84.62 (0.28)	79.68 (2.55)	71.99 (1.24)	83.78 (0.28)	84.86 (0.22)	85.03 (0.61)

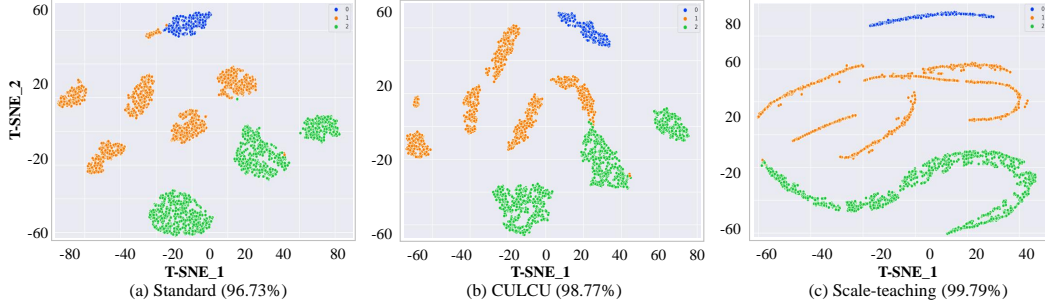


Figure 6: t-SNE visualization of the learned embeddings on the FD-A dataset with Sym 50% noisy labels (values in parentheses are the test classification accuracies).

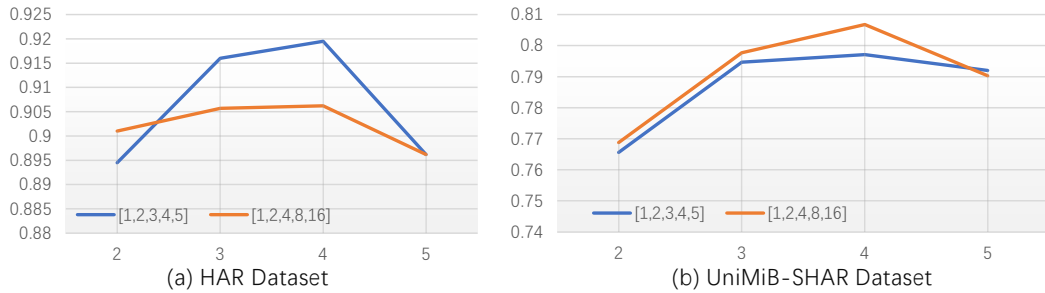


Figure 7: Multi-scale sampling strategies analysis under Sym 50% noisy labels.

time-series label-noise learning can make the embeddings of different classes more discriminative, thus facilitating clean sample selection and noisy label correction.

Impact of downsampling scale sequence list. Scale-teaching can be performed using a variety of different downsampling scales for label-noise learning. Based on the experience of [19, 44] on time series classification and prediction tasks, we utilize the downsampling scales of [1,2,4] for the experimental analyses of Scale-teaching. However, for real-world scenarios that actually contain noisy labels, it is generally not possible to perform hyperparametric analyses using a clean-labeled validation set. To facilitate the analysis, in this paper, we use the classification performance of the test set for multi-scale hyperparameter analyses. However, to avoid test set information leakage, we do not use the hyperparameter analysis result for Scale-teaching in our experiments. We use two multi-scale sampling strategies for analyses, which are (1) {[1,2], [1,2,3], [1,2,3,4], [1,2,3,4,5]}; (2) {[1,3], [1,2,4], [1,2,4,8], [1,2,4,8,16]}. From Figure 7, we find that Scale-teaching using four different scales for training has the highest classification accuracy, which indicates that more input scales do not necessarily make the classification performance better. In addition, using three or four scales of sequences can effectively improve the classification performance of Scale-teaching compared with using two different scales.

Impact of input scales of sequences order. Scale-teaching employs a finer-to-coarser strategy for cross-scale embedding fusion. Intuitively, when a single scale is used for classification, the original single scale (finer) time series is better overall because it does not discard the original sequence information compared to coarser scale time series. Therefore, Scale-teaching is trained using the finer-to-coarser cross-scale fusion strategy. To analyze the difference in classification performance between different fusion directions, we subtract the classification accuracy using the finer-to-coarser and coarser-to-finer training approaches, and the specific results are shown in Figure 8. We can find that the classification performance of finer-to-coarser is better overall, which is due to its ability to use a single fine-scale sequence with an excellent classification performance from the beginning to gradually promote the classification performance of multiscale fusion embeddings.

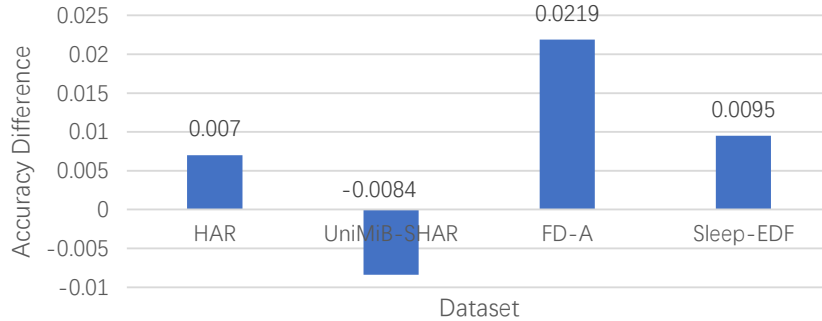


Figure 8: The cross fusion direction of input scale series analysis under Sym 50% noisy labels.

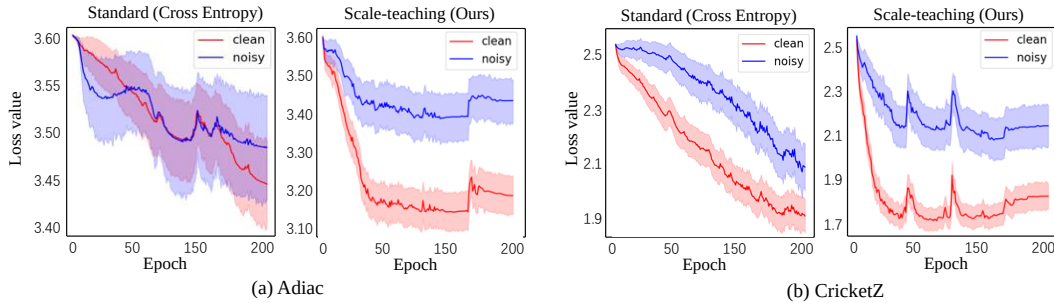


Figure 9: The change of loss values for clean and noisy time series samples under Aysm 40% noise labels. The solid line and shading indicate the mean and standard deviation loss values of all clean (or noisy) training samples within each epoch.

F Small-loss Visualization

The small-loss criterion has been extensively validated for clean label selection in label-noise learning for computer vision. To further analyze the application of the small-loss criterion in time series data, we provide the change of loss values of the models trained by the Standard method and Scale-teaching on Adiac and CricketZ UCR datasets (as in Figure 9). Also, we visualize the probability distributions of the ground-truth clean and noisy (corrupted) sample loss values on the test set with different training strategies. Specifically, Figures 10 and 11 show the loss probability distributions of the models trained by different strategies on the HAR dataset and UniMiB-SHAR with Aysm 40% noisy labels. Both red (clean) and blue (corrupted) in Figure 10 and Figure 11 contain two peaks, which indicate that some correctly labeled samples are still difficult to learn (large loss) and some incorrectly labeled samples are also easy to learn (small loss). Compared with the Standard method (Figure 10 (a) and Figure 11 (a)), Scale-teaching (Figure 10 (b) and Figure 11 (b)) can clearly distinguish clean and noisy samples by the loss value distribution, further validating the robustness of the multi-scale embeddings to cope with time-series noisy labels.

G Other Analysis

The test F1-score results of ablation study. Following [4], we select the averaged F1-score on the test set as a new metric for ablation analysis in Section 4.5. Hence, we give the corresponding test classification F1-score (%) in Tables 7 and 8.

Running time analysis. We select two datasets for running the time-consuming analysis, the FD-A dataset with the largest sequence length and the Sleep-EDF dataset with the largest samples. We performed the running time statistics on the NVIDIA GeForce RTX 3090 GPU using all baselines, and the results are shown in Table 9. On the FD-A dataset with the longest sequence length, Co-teaching and CULCU take essentially twice as long to run as the Standard method because they use

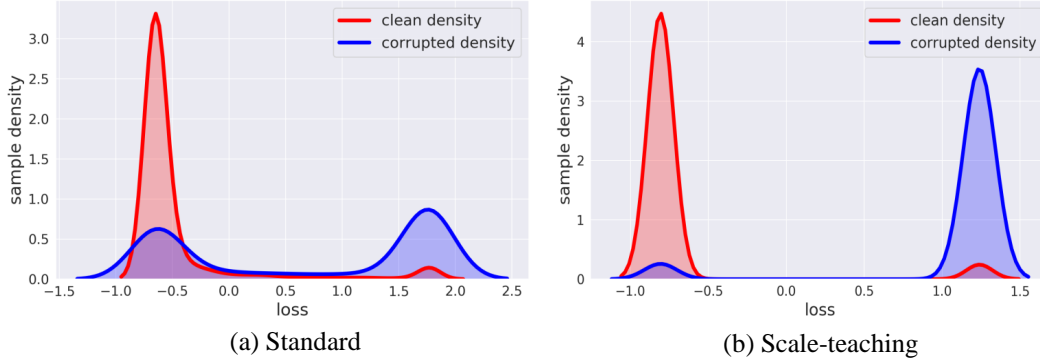


Figure 10: The loss value probability distributions visualization on HAR dataset with Asym 40% noisy labels.

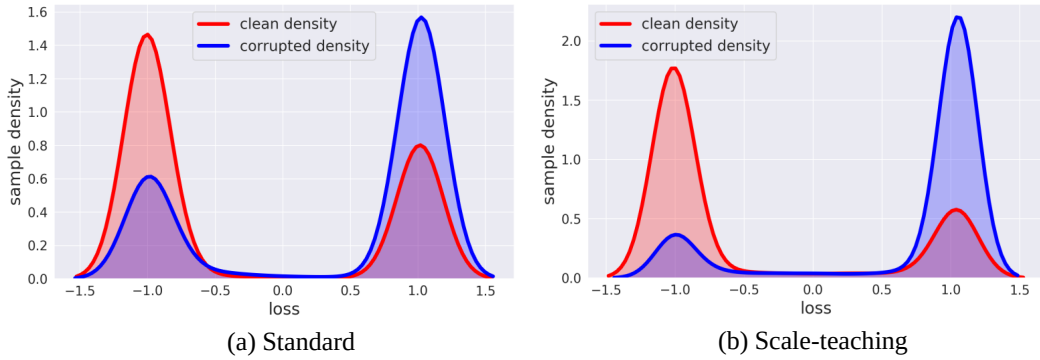


Figure 11: The loss value probability distributions visualization on UniMiB-SHAR dataset with Asym 40% noisy labels.

Table 7: The test classification F1-score (%) results of ablation study (values in parentheses denote drop F1-score).

Method	HAR		UniMiB-SHAR	
	Sym 50%	Asym 40%	Sym 50%	Asym 40%
Scale-teaching	90.05	89.14	77.56	65.89
w/o cross-scale fusion	88.16 (-1.89)	87.05 (-2.09)	68.23 (-9.33)	57.76 (-8.13)
only single scale	87.56 (-2.49)	86.75 (-2.39)	66.87 (-10.69)	54.12 (-11.77)
w/o graph learning	87.79 (-2.26)	87.41 (-1.73)	74.62 (-2.94)	63.15 (-2.74)
w/o moment	89.34 (-0.71)	88.27 (-0.87)	76.67 (-0.89)	64.92 (-0.97)
w/o dynamic threshold	88.93 (-1.12)	88.29 (-0.85)	73.11 (-4.45)	64.76 (-1.17)

Table 8: The test classification F1-score (%) results on four individual large datasets without noisy labels. The best results are **bold**, and the second best results are underlined.

Dataset	Standard	Mixup	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
HAR	93.27	95.39	93.75	93.19	92.91	93.71	<u>94.72</u>	94.18
UniMiB-SHAR	86.37	80.17	84.43	84.03	66.54	<u>89.19</u>	<u>86.45</u>	93.62
FD-A	99.93	99.91	99.96	64.05	90.14	<u>99.82</u>	<u>99.95</u>	99.96
Sleep-EDF	81.99	82.11	82.52	83.07	77.67	82.17	<u>83.26</u>	84.76

Table 9: Training time (hours) analysis using the FD-A and Sleep-EDF datasets with Asym 40% noisy labels.

Dataset	Standard	Mixup	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching		
								[1,2,4]	[1,4,16]	[1,8,32]
FD-A	0.37	0.42	0.79	0.63	0.90	0.42	0.87	1.06	0.86	0.82
Sleep-EDF	0.54	0.83	1.09	2.04	1.64	0.73	1.47	2.02	1.76	1.60

Table 10: Comparison with classification methods without label noise learning strategy. The best test classification accuracy (%) results are **bold**, and the second best results are underlined.

Dataset	HAR					FD-A				
	0	Sym 20%	Sym 50%	Asym 40%	Ins 40%	0	Sym 20%	Sym 50%	Asym 40%	Ins 40%
Boss [15]	72.34	62.55	56.11	53.29	52.34	69.75	64.75	57.95	61.99	62.25
Rocket [60]	95.29	<u>92.93</u>	90.04	<u>82.53</u>	<u>90.43</u>	99.99	99.71	97.01	89.75	97.98
FCN [57]	93.74	92.13	83.99	75.59	83.56	99.56	98.89	96.63	96.12	99.36
Scale-teaching	<u>94.72</u>	93.93	90.17	89.62	91.58	<u>99.98</u>	99.93	99.38	99.55	99.82

two encoders. Furthermore, although SREA uses a single network training, it utilizes a decoder for the unsupervised reconstruction task of the original time series, which significantly increases training time on the FD-A dataset with longer sequences. running time is higher than Co-teaching. The Scale-teaching paradigm uses multiple encoders for training and has an additional noisy label correction module, which is expected to increase the training time. Nevertheless, the larger the sampling scale (coarse scale) of the training data used by the Scale-teaching paradigm, the lower the training elapsed time of its model. For example, with input scales of [1, 8, 32], the training time of Scale-teaching is lower than that of CULCU and SREA. On the SleepEEG dataset with the largest number of samples, we find that FINE with an encoder has a higher running time because FINE using all training samples to select clean labels is time-consuming when the sample size is large. In contrast, the runtime of Scale-teaching is lower than FINE. Also, when the input scales are set to [1,8,32], the runtime of Scale-teaching is lower than SREA.

It is worth noting that when Scale-teaching is trained using two scales, such as [1,2] or [1,16], its training run time decreases further. From the analysis in Appendix E, it is clear that Scale-teaching using three different scales generally performs better than two scales for classification with noisy labels. In addition, the classification performance of [1,2,4], [1,4,16], and [1,8,32] when Scale-teaching is trained using three different scales has less difference in classification performance on datasets with longer sequences (e.g., FD-A and Sleep-EDF). The above results indicate that the Scale-teaching paradigm has a greater advantage in runtime on time-series datasets with longer sequences.

Robustness analysis. Three time-series supervised classification methods (Boss [15], Rocket [60] and FCN [57]) and the Scale-teaching paradigm are chosen for robustness analysis against time-series noise labels. Boss [15] is a time series classification method based on similarity search, which can effectively mitigate the negative impact of noise (e.g., adding Gaussian noise) in time series values on classification. Rocket [60] uses a large number of randomly initialized convolution kernels to extract time series features, and employs the extracted features to classify time series using a machine learning classifier (e.g., Ridge classifier). FCN [57] is the encoder used by Scale-teaching, which is a time series classification method based on DNNs. As shown in Table 10, the classification performance of the Scale-teaching paradigm using FCN as encoders is better than that of Boss, Rocket and FCN in the presence of noisy labels. It is worth noting that both Boss and Rocket training processes are independent of the optimization of DNNs. However, their classification performance is still reduced due to the influence of noisy labels. In addition, the encoder of the Scale-teaching paradigm can be designed flexibly, such as using ResNet [1], InceptionTime [65] and OS-CNN [20] in the field of time series classification. In other words, using better robustness encoders, the classification performance of Scale-teaching can be further improved with time-series noise labels.

Table 15: The detailed test classification accuracy (%) results on UEA 30 archive with Sym 20% noisy labels.

ID	Dataset	Standard	MixUp	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
1	ArticularyWordRecognition	89.33	81.04	81.77	41.07	90.13	87.87	84.33	92.73
2	AtrialFibrillation	26.67	25.67	26.67	33.33	26.67	26.67	22.67	28.00
3	BasicMotions	97.10	97.40	99.25	91.00	97.50	98.00	97.50	95.50
4	CharacterTrajectories	98.16	98.49	92.59	7.45	96.54	98.93	98.66	98.93
5	Cricket	94.72	94.17	91.53	41.11	97.50	96.67	93.89	97.28
6	DuckDuckGeese	48.40	52.48	48.60	24.40	46.80	50.00	49.60	45.60
7	EigenWorms	53.56	55.42	62.88	34.66	61.80	64.58	64.10	61.25
8	Epilepsy	84.96	80.28	93.33	82.46	88.70	86.38	93.33	87.59
9	EthanolConcentration	24.15	25.48	26.43	24.71	25.11	24.87	24.39	27.42
10	ERing	73.11	73.70	70.04	16.67	64.59	72.59	61.70	72.77
11	FaceDetection	52.90	51.93	52.36	51.96	50.93	52.30	52.69	53.09
12	FingerMovements	50.80	50.44	52.76	52.20	50.20	52.80	52.30	54.68
13	HandMovementDirection	28.92	32.22	35.14	24.05	19.46	28.11	29.16	35.35
14	Handwriting	35.71	28.55	30.71	25.31	16.63	28.19	34.12	41.35
15	Heartbeat	52.00	52.93	57.46	72.10	62.75	50.17	66.39	51.76
16	InsectWingbeat	62.31	53.90	64.74	63.92	51.56	64.63	64.75	63.85
17	JapaneseVowels	88.99	87.61	94.70	49.73	97.41	93.51	97.19	95.62
18	Libras	74.60	70.07	77.00	10.78	73.22	76.78	70.83	79.67
19	LSST	49.68	48.91	50.23	51.61	35.53	50.48	51.91	55.16
20	MotorImagery	51.00	50.88	51.60	50.80	55.92	52.80	50.10	50.84
21	NATOPS	80.78	80.40	90.56	25.78	89.33	82.67	87.17	91.11
22	PenDigits	96.04	97.58	97.33	98.30	97.90	98.27	98.16	98.12
23	PEMS-SF	62.08	62.87	60.06	14.91	64.97	62.08	61.10	63.05
24	PhonemeSpectra	21.89	23.46	22.25	18.00	6.23	22.82	23.74	26.50
25	RacketSports	76.84	73.50	78.64	25.00	75.74	77.24	79.41	76.87
26	SelfRegulationSCP1	76.81	77.49	78.91	56.86	49.83	78.43	77.55	81.28
27	SelfRegulationSCP2	47.47	49.24	48.32	50.89	50.00	47.44	50.66	51.78
28	SpokenArabicDigits	96.17	96.21	98.83	27.79	98.52	98.67	98.71	99.18
29	StandWalkJump	44.00	40.00	36.00	33.33	45.33	42.67	33.40	34.13
30	UWaveGestureLibrary	76.90	76.97	74.61	12.50	74.29	77.75	66.97	78.12
Avg Acc		63.87	62.98	64.84	40.42	62.04	64.81	64.55	66.29
Avg Rank		5.03	5.20	3.83	6.37	4.77	3.73	4.00	2.73
P-value		6.61E-04	3.33E-04	2.69E-02	2.37E-05	1.14E-02	2.63E-02	3.93E-02	-

Table 16: The detailed test classification accuracy (%) results on UEA 30 archive with Sym 50% noisy labels.

ID	Dataset	Standard	MixUp	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
1	ArticularyWordRecognition	66.64	55.57	52.12	7.53	67.33	57.40	49.63	68.68
2	AtrialFibrillation	28.00	26.67	29.33	33.33	29.33	30.67	32.00	36.00
3	BasicMotions	56.00	58.20	70.75	54.50	81.00	59.00	71.25	60.20
4	CharacterTrajectories	90.62	87.56	67.75	6.73	95.28	93.82	97.03	97.23
5	Cricket	67.44	60.00	68.89	8.33	84.28	66.67	68.06	87.27
6	DuckDuckGeese	32.00	36.64	29.60	26.00	34.00	33.20	33.20	33.84
7	EigenWorms	49.62	45.37	55.65	32.98	38.78	56.79	56.09	49.28
8	Epilepsy	57.39	57.25	64.16	36.38	62.62	60.00	63.55	64.36
9	EthanolConcentration	25.20	25.32	26.05	25.02	27.56	25.62	24.58	26.27
10	ERing	44.30	44.39	40.63	16.67	39.35	44.52	38.54	45.10
11	FaceDetection	48.54	49.81	49.30	49.51	48.94	49.11	48.58	50.72
12	FingerMovements	50.40	50.20	51.80	51.40	50.60	50.80	51.30	54.40
13	HandMovementDirection	26.22	25.51	25.41	26.11	23.51	26.57	24.35	27.08
14	Handwriting	19.54	19.53	19.99	13.41	7.53	19.18	18.28	21.06
15	Heartbeat	55.32	53.52	52.20	54.44	55.40	54.63	53.27	48.62
16	InsectWingbeat	49.30	32.34	58.02	52.25	31.07	53.97	58.13	52.01
17	JapaneseVowels	60.14	59.28	73.97	15.03	70.65	66.43	78.11	79.23
18	Libras	47.09	43.98	51.39	6.67	44.78	50.11	40.06	49.78
19	LSST	47.29	44.58	46.21	47.89	34.35	46.33	47.92	48.75
20	MotorImagery	50.84	51.48	50.10	51.40	52.00	50.60	49.70	49.80
21	NATOPS	54.33	52.53	60.17	16.67	59.80	53.44	58.06	59.56
22	PenDigits	93.38	85.29	95.92	92.85	92.74	96.80	96.53	93.89
23	PEMS-SF	41.20	40.55	32.96	14.45	42.43	43.24	37.57	41.27
24	PhonemeSpectra	19.08	19.11	19.94	11.48	3.92	19.69	19.23	20.09
25	RacketSports	52.50	51.03	52.80	24.08	53.29	52.89	56.58	54.21
26	SelfRegulationSCP1	47.41	42.68	48.86	48.60	49.97	48.33	57.93	58.08
27	SelfRegulationSCP2	48.13	48.22	47.61	48.78	50.00	49.22	50.00	48.73
28	SpokenArabicDigits	85.64	69.95	96.55	96.16	99.23	95.66	97.59	97.69
29	StandWalkJump	38.67	37.87	40.67	33.33	42.67	42.67	37.33	44.00
30	UWaveGestureLibrary	50.41	48.60	45.45	12.50	53.91	49.00	37.94	52.61
Avg Acc		50.09	47.43	50.81	33.82	50.88	51.55	51.75	53.99
Avg Rank		5.17	5.73	4.23	6.23	3.93	3.83	4.30	2.43
P-value		2.98E-04	7.40E-05	1.59E-02	9.35E-05	1.67E-02	1.08E-02	3.75E-02	-

Table 17: The detailed test classification accuracy (%) results on UEA 30 archive with Asym 40% noisy labels.

ID	Dataset	Standard	MixUp	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
1	ArticularyWordRecognition	66.33	62.40	55.87	17.27	69.40	63.67	53.73	70.44
2	AtrialFibrillation	21.67	33.33	33.33	33.33	11.33	33.67	32.67	34.67
3	BasicMotions	66.00	62.30	67.25	49.50	69.00	64.00	61.75	65.10
4	CharacterTrajectories	61.08	60.01	57.42	19.05	87.78	64.29	61.35	88.34
5	Cricket	72.22	71.56	70.97	50.00	73.44	72.78	68.61	80.56
6	DuckDuckGeese	43.20	42.96	44.80	24.00	43.68	45.20	44.60	38.24
7	EigenWorms	41.75	34.75	51.34	37.86	43.56	41.68	50.38	42.47
8	Epilepsy	62.32	63.01	63.48	47.25	61.01	61.45	64.71	58.70
9	EthanolConcentration	23.04	23.85	23.61	25.02	24.78	24.33	25.57	27.70
10	ERing	60.30	60.37	42.74	39.47	45.56	59.11	43.96	61.74
11	FaceDetection	49.88	50.64	51.12	51.06	50.32	51.07	50.12	51.61
12	FingerMovements	47.76	49.92	48.50	49.80	49.00	48.20	50.19	50.96
13	HandMovementDirection	28.97	31.24	30.41	28.38	31.46	29.19	29.73	29.03
14	Handwriting	21.61	23.92	24.74	21.03	10.67	22.99	25.69	26.98
15	Heartbeat	55.22	57.46	55.61	72.20	61.52	55.12	55.51	56.68
16	InsectWingbeat	43.40	38.07	45.34	46.32	48.78	47.81	50.34	51.87
17	JapaneseVowels	61.62	58.63	62.46	36.81	65.97	64.27	73.76	70.02
18	Libras	57.47	57.00	53.39	8.67	54.33	59.44	45.72	63.22
19	LSST	42.11	42.77	41.70	43.67	32.79	43.67	43.74	29.10
20	MotorImagery	48.80	50.24	51.32	53.00	52.60	49.60	49.70	53.20
21	NATOPS	57.00	55.29	58.65	16.67	64.89	55.89	63.22	65.13
22	PenDigits	78.76	67.36	92.78	84.05	91.07	89.18	92.23	93.57
23	PEMS-SF	50.20	51.38	42.60	14.45	50.87	50.76	47.86	51.45
24	PhonemeSpectra	17.71	19.05	18.70	14.41	5.11	18.02	18.65	19.52
25	RacketSports	57.50	59.13	54.30	27.50	55.26	58.16	56.32	54.21
26	SelfRegulationSCP1	63.47	66.21	64.94	60.96	49.83	66.42	68.10	66.30
27	SelfRegulationSCP2	49.04	51.24	51.20	52.16	50.00	52.11	51.26	51.22
28	SpokenArabicDigits	64.19	60.64	79.13	72.11	99.04	79.42	88.10	93.85
29	StandWalkJump	38.67	34.67	39.33	33.33	33.33	40.00	39.33	36.27
30	UWaveGestureLibrary	53.56	53.36	55.38	12.50	57.81	53.69	45.53	55.84
	Avg Acc	50.16	49.76	51.08	38.06	51.47	52.17	51.75	54.60
	Avg Rank	5.60	4.77	4.40	6.13	4.20	4.00	3.97	2.73
	P-value	3.81E-03	6.17E-03	1.63E-02	9.33E-05	1.36E-02	2.62E-02	3.88E-02	-

Table 18: The detailed test classification accuracy (%) results on UEA 30 archive with Ins 40% noisy labels.

ID	Dataset	Standard	MixUp	Co-teaching	FINE	SREA	SELC	CULCU	Scale-teaching
1	ArticularyWordRecognition	67.27	57.39	60.93	9.20	68.67	61.73	57.10	75.40
2	AtrialFibrillation	28.00	29.33	30.00	33.33	32.00	28.00	34.67	32.00
3	BasicMotions	77.00	73.00	81.75	39.00	80.90	77.50	74.75	78.50
4	CharacterTrajectories	82.52	69.46	66.33	5.22	85.38	81.92	83.48	87.47
5	Cricket	80.28	78.28	79.31	26.11	92.50	81.11	76.81	92.78
6	DuckDuckGeese	38.80	41.60	36.40	20.00	39.20	39.60	42.00	37.20
7	EigenWorms	33.44	54.63	60.38	32.67	43.66	43.21	60.31	54.81
8	Epilepsy	65.80	64.55	72.03	67.10	73.04	71.74	78.64	79.36
9	EthanolConcentration	24.82	26.40	27.75	24.71	25.17	26.69	26.81	28.37
10	ERing	53.69	52.48	48.48	36.79	49.82	56.44	44.59	54.15
11	FaceDetection	50.05	50.04	50.80	50.57	49.81	50.22	50.93	51.00
12	FingerMovements	51.40	51.60	51.10	51.20	49.60	51.40	48.60	49.60
13	HandMovementDirection	25.31	30.11	26.49	19.73	19.73	28.92	25.95	31.62
14	Handwriting	21.35	22.54	22.92	17.93	6.17	21.51	23.19	23.01
15	Heartbeat	56.39	58.01	66.63	57.69	58.50	56.20	59.02	51.00
16	InsectWingbeat	47.94	39.18	57.07	55.93	36.95	57.12	59.40	58.32
17	JapaneseVowels	68.05	66.46	65.36	27.89	77.03	73.35	78.97	81.54
18	Libras	46.04	51.56	47.00	9.33	48.67	49.44	41.11	50.18
19	LSST	48.16	47.72	46.43	49.04	33.58	48.78	49.11	50.47
20	MotorImagery	49.20	51.00	49.90	50.00	51.00	52.00	51.90	47.80
21	NATOPS	57.78	56.80	58.78	26.67	67.24	57.67	58.34	59.98
22	PenDigits	81.99	70.68	93.43	91.81	92.59	91.18	93.29	96.69
23	PEMS-SF	42.89	43.86	35.14	16.42	47.86	43.82	41.49	44.35
24	PhonemeSpectra	16.50	15.82	17.37	8.40	3.22	17.15	17.48	17.42
25	RacketSports	59.08	59.55	64.68	27.50	60.11	58.55	58.62	59.89
26	SelfRegulationSCP1	60.63	60.38	56.02	48.12	49.90	65.80	66.00	54.54
27	SelfRegulationSCP2	51.69	50.87	50.63	51.11	50.00	51.89	50.22	52.11
28	SpokenArabicDigits	74.85	67.24	89.70	83.37	98.76	86.68	87.87	97.53
29	StandWalkJump	40.00	42.67	39.33	32.00	42.67	38.67	42.00	40.00
30	UWaveGestureLibrary	67.19	67.90	60.56	12.50	67.74	67.69	55.69	69.96
	Avg Acc	52.27	51.70	53.76	36.04	53.38	54.53	54.61	56.90
	Avg Rank	5.20	4.77	4.33	6.60	4.27	4.20	3.77	2.60
	P-value	6.08E-04	2.92E-03	1.20E-02	2.55E-05	5.52E-03	1.08E-02	3.47E-02	-

