

# Deep Contrastive One-Class Time Series Anomaly Detection\*

Rui Wang<sup>†</sup>   Chongwei Liu<sup>†</sup>   Xudong Mou<sup>†</sup>   Kai Gao<sup>‡</sup>   Xiaohui Guo<sup>§</sup>  
 Pin Liu<sup>¶</sup>   Tianyu Wo<sup>||</sup>   Xudong Liu<sup>†§</sup>

## Abstract

The accumulation of time-series data and the absence of labels make time-series Anomaly Detection (AD) a self-supervised deep learning task. Single-normality-assumption-based methods, which reveal only a certain aspect of the whole normality, are incapable of tasks involved with a large number of anomalies. Specifically, Contrastive Learning (CL) methods distance negative pairs, many of which consist of both normal samples, thus reducing the AD performance. Existing multi-normality-assumption-based methods are usually two-staged, firstly pre-training through certain tasks whose target may differ from AD, limiting their performance. To overcome the shortcomings, a deep **Contrastive One-Class Anomaly** detection method of time series (COCA) is proposed by authors, following the normality assumptions of CL and one-class classification. It treats the original and reconstructed representations as the positive pair of negative-sample-free CL, namely “sequence contrast”. Next, invariance terms and variance terms compose a contrastive one-class loss function in which the loss of the assumptions is optimized by invariance terms simultaneously and the “hypersphere collapse” is prevented by variance terms. In addition, extensive experiments on two real-world time-series datasets show the superior performance of the proposed method achieves state-of-the-art.

## 1 Introduction

Within cyber-physical systems, sensor-equipped devices generate time-series data that contains massive status information, making it possible to detect unexpected er-

rors and reduce maintenance costs in data-driven ways. Anomaly Detection (AD) plays an increasingly important role in this context, which refers to detecting instances that are significantly dissimilar to the majority [13]. Though the performance of deep learning methods is superior to shallow ones [21], labeling the outlier from quantities of temporal data could be costly and tricky. So, AD is usually considered an unsupervised learning problem in which learning representation for discerning anomalies relies on some normality assumptions. For example, autoencoder-based [19] methods assume normal samples are better restructured from the latent space than abnormal ones. Similarly, one-class classification methods [24] assume that the normal samples come from a single (abstract) class that could accurately describe the so-called “normality”. However, these normality assumptions may be one-sided, some of which are just inspired by the pretext task of self-supervised representation learning. Meanwhile, there are various time-series anomalies including point anomalies (global or local), subsequence anomalies, and anomaly time series [4] (Fig. 3), thus it is not sufficient to detect all based on one normality assumption alone.

In particular, contrastive learning-based AD methods are emerging. [10] directly treats the InfoNCE loss of CPC [20] as the anomaly score for image AD, contrasting the context vector with the future representation vector. NeuTraL AD [22] devises a contrastive loss specific to a fixed set of learnable transformations and regards the training loss as the anomaly score, contrasting the transformed samples (views) with the original ones in the representation space. The single-assumption-based CL AD methods above assume that more mutual information exists between normal comparison objects than anomalous ones. However, pairs transformed from different normal samples are treated as negative ones, pushing away many normal samples inside and not capturing shared information in the same class, similar to [5]. It goes against the very nature of AD, i.e., extracting features common to the vast majority of normal samples, thus leading to a decline in AD performance.

Indeed some scholars combine these normality as-

\*The full version of the paper can be accessed at <https://arxiv.org/abs/2207.01472>

<sup>†</sup>School of Computer Science and Engineering, Beihang University, Beijing, China. {ruiking, liucw, mxid}@buaa.edu.cn, liuxd@act.buaa.edu.cn

<sup>‡</sup>University of New South Wales, Sydney, Australia. kai.gao@unsw.edu.au

<sup>§</sup>Hangzhou Innovation Institute, Beihang University, Hangzhou, China. Xiaohui Guo is the corresponding author. guoxh@buaa.edu.cn

<sup>¶</sup>School of Information Engineering, China University of Geosciences in Beijing, Beijing, China. liupin@cugb.edu.cn

<sup>||</sup>College of Software, Beihang University, Beijing, China. woty@act.buaa.edu.cn

sumptions into some compound ones to learn more expressive representations for downstream AD tasks. For instance, Deep SVDD [24] realizes a deep one-class classification framework for AD with deep features or representations learned by a pre-trained autoencoder. [27] presents the two-stage one-class classifier on contrastive representations and points out a subtle but important observation, i.e., the uniformity property of contrastive representation may hurt the one-class AD performance. Even though, we argue that learning representation is distinct from capturing the normality and anomalies' underlying data regularities, formally they are two discrepant optimization objectives. Therefore, with representation learning and outlier discriminating separated, the two-stage AD methods' performance is limited. In addition, these AD methods are originally proposed in the computer vision domain, lacking temporal dependencies, thus generalizing them simply into time-series AD tasks is meaningless.

To address the above issues, we propose a one-stage negative-sample-free deep **Contrastive One-Class Anomaly** (COCA) detection model for time-series data. As shown in Fig. 1, first, the original training data is augmented, making it easier to isolate anomalies from normal samples. Next, the augmented time series is encoded through a multi-layer temporal convolution neural network and then put into a Seq2Seq model in the latent space to learn the critical characteristics of time series, i.e., temporal dependencies. The key to CL is to pull contrasting objects (positive pairs) closer in the representation space, and researchers use a variety of positive pairs, such as context/future [20], different augmentations [6], and context/mask [1]. Here, we regard the representation in the latent space and the representation reconstructed by the Seq2Seq model as positive pairs and name it "sequence contrast". Note that it's different from an autoencoder, as the latter is a generative method, which performs the reconstruction of original data or so-called pixel-level generation [6], carrying massive unnecessary details to downstream tasks. Finally, the positive pairs are fed to a learnable nonlinear projection layer to obtain their projections respectively.

The model is trained via a contrastive one-class loss function with two terms: *invariance* and *variance*. The invariance term is to maximize the cosine similarity between the one-class center, latent representations, and seq2seq outputs, instead of adjusting the hyperparameters to balance the loss contribution of one-class and contrastive learning as in most multi-task learning. The variance term is borrowed from [2], and the variance of the within-batch representations is maintained above a given threshold by a hinge loss to

avoid "hypersphere collapse" without negative sample pairs, which also solves the difficulty of identifying negative pairs in AD. In practice, the invariance term is treated as the anomaly score for AD. In conclusion, COCA combines the two normality assumptions that latent and reconstructed representations 1) have greater mutual information and 2) belong to a single class, without pre-training. We summarize our contributions as follows:

- A novel normality assumption that combines CL and one-class classification for time-series AD.
- A new time-series CL paradigm namely "sequence contrast". By analyzing the problems solved with CL, we clarify that its essence is the representation, rather than the compared pairs or the negative examples.
- A novel contrastive one-class loss function to optimize both contrastive learning and one-class classification, and avoid "hypersphere collapse" at the same time.
- Extensive experiments performed on two datasets show that the proposed COCA leads to a new state-of-the-art in time-series AD.

## 2 Related Work

This section contains a brief introduction of recent works in contrastive learning and deep anomaly detection.

**Contrastive Learning.** The recent renaissance of contrastive learning began with CPC [20] proposed InfoNCE, which pulls positive samples closer and distances negative samples, though relying on a large number of negative samples to learn a good representation. [28] summarizes two key properties of contrastive learning: 1) alignment: similar samples have similar representations (pull positive pair) and 2) uniformity: representations follow a uniform distribution on the hypersphere (push negative pair). On the one hand, BYOL [12], SwAV [5], and SimSiam [7] achieve uniformity in contrastive learning without using negative samples. On the other hand, SimCLR [6] and TS-TCC [11] align augmented data representations to learn invariant representations for visual data and time series, respectively. Also, TS-TCC uses a temporal contrasting module to address the temporal dependencies of time series. Although all these contrastive learning approaches have successfully improved representation learning for visual data and time series, they could be inapplicable to time-series AD. For example, contradictions exist between the uniformity of contrastive learning and the class imbalance of anomaly detection.

**Deep Anomaly Detection.** Recently, deep learning for anomaly detection has been regarded as a new research frontier of the AD field. Deep anomaly detection methods can roughly be divided into two categories: deep learning for feature extraction and learning feature representations of normality [21]. Deep learning for feature extraction is a two-staged learning method that uses deep methods to learn representations for downstream anomaly detection. However, it does not directly address the anomaly detection task, so the representations learned in the pre-training may be detrimental to anomaly detection. Learning feature representations of normality couples representations learning with anomaly scoring in some way, such as GANs-based [25], autoencoder-based [19], one-class classification-based [24], clustering-based [30], saliency map-based [24], and contrastive learning-based [10, 22] methods. The key to these methods lies in the assumption of normality/anomaly, and some assumptions of normality are inspired by the pretext task of self-supervised learning. For instance, GANs-based methods assume normal samples are better generated from the latent space of the generative network than anomalies. However, the normal sample assumption of these methods may explain only one aspect of overall normality, respectively. Uniquely, COCA does not resort to pre-training and organically integrates the normality assumption of one-class classification and contrastive learning to detect anomalies for time-series data.

### 3 Methodology

This section describes the proposed COCA in detail, including the structure, objective, and its relation to contrastive learning.

**3.1 Problem Definition.** Given a set of time series  $\mathcal{D} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$ ,  $\mathbf{X}_i = \{x_1, x_2, \dots, x_T\}$  is a time series of length  $T$ , where  $x_j \in \mathbb{R}^d$  is a  $d$ -dimensional vector. Since sliding windows are generally used to divide time series into length- $T$  sequences,  $T$  has been called the sliding window length, as well.  $d = 1$  means that the time series is univariate, and  $d > 1$  for multivariate. In time-series AD, the anomaly score  $\mathbf{S}_i$  of  $\mathbf{X}_i$  is calculated by the AD model such that the higher  $\mathbf{S}_i$  is, the more likely it is an anomalous time series.

**3.2 Architecture.** Fig. 1 shows the architecture of the COCA model. The time series  $\mathbf{X}_i$  from an augmented training set of the raw dataset is passed to a multi-layer temporal convolution feature encoder  $f_\theta : \mathcal{X} \mapsto \mathcal{Z}$  which takes as input time series  $\mathbf{X}_i$  of length  $T$  and outputs latent representations  $z_1, \dots, z_L$  for  $L$  time-steps, potentially with a lower temporal resolution,

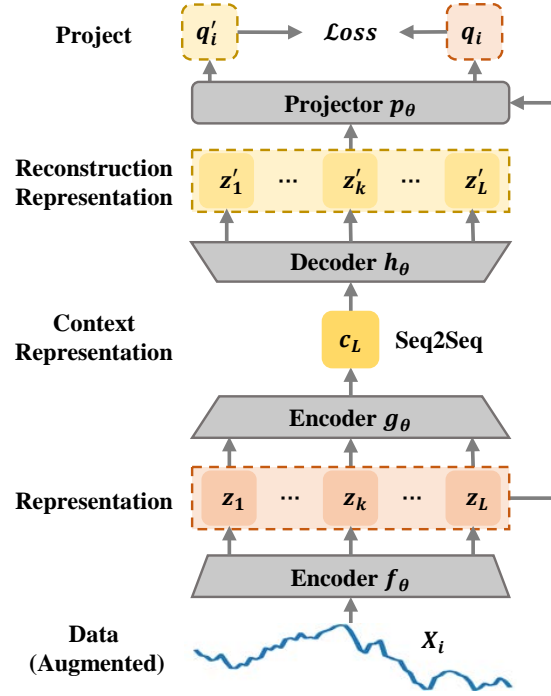


Figure 1: The overall architecture of the proposed COCA model.

i.e.  $T > L$ . They are then fed to a Seq2Seq encoder  $g_\theta : \mathcal{Z} \mapsto \mathcal{C}$  to summarize all  $z_{\leq L}$  as context vectors  $c_L$  and then a Seq2Seq decoder  $h_\theta : \mathcal{C} \mapsto \mathcal{Z}$  produces reconstruction representations  $z'_1, \dots, z'_L$  for  $L$  time-steps to learn temporal dependencies. Furthermore, latent representations  $z_k$  and seq2seq outputs  $z'_k$  are fed to a learnable nonlinear projector  $p_\theta : \mathcal{Z} \mapsto \mathcal{Q}$  to output projections  $q$  and  $q'$ . The output of the projector is used to calculate the loss (see next sub-section 3.3) to maximize the similarity between  $q$  and  $q'$  concerning the one-class center  $Ce \in \mathcal{Q}$  to combine the two normality assumptions: contrastive learning-based and one-class classification-based.

**Time-Series Augmentation.** Data augmentation helps improve the performance of AD methods because it not only increases the volume of train data but also makes it easier to isolate anomalies [27]. In this paper, jittering (noise addition) and scaling (pattern-wise magnitude change) are applied to expand the training set. Notably, the jittering and scaling hyper-parameters should be carefully chosen according to the nature of the time-series anomalies.

**Feature encoder.** The encoder network has a 2-block temporal convolutional architecture, each block comprising a Conv1D layer, a BatchNormalization (BN)

layer, a ReLU activation function, and a MaxPool1D layer, where the first block also contains a Dropout layer. The time series input to the encoder should be normalized to zero mean and unit variance.

**Seq2Seq.** The Seq2Seq consists of an encoder and a decoder. The encoder is a 3-layer Long Short-Term Memory (LSTM) and the decoder is a 3-layer LSTM followed by a fully-connected (FC) layer. In this paper, the hidden space representation length  $L < 20$ , therefore LSTM can meet the needs of the context representation, while for long sequences, more recent advancements in Seq2Seq modeling such as self-attention networks or the Transformer model could help improve results further.

**Projector.** The projector uses an MLP with one hidden layer applied BN and ReLU to map representations to the space where contrastive one-class loss is calculated.

**3.3 The COCA Objective.** The COCA objective consists of *invariance* and *variance* terms. The invariance term is to maximize the cosine similarity between the one-class center  $Ce$ , representations  $q_i$ , and seq2seq outputs  $q'_i$  in the projection space  $\mathcal{Q}$ , and the variance term avoids “hypersphere collapse” without negative sample pairs.

Before explaining the invariance term of the COCA objective, it is necessary to state the optimization objectives of one-class classification and contrastive learning without negative pairs.

**One-class classification.** The optimization objective of Deep SVDD [24], a representative method for one-class classification, is defined as:

$$(3.1) \quad \mathcal{L}_{svdd} = \frac{1}{N} \sum_{i=1}^N \|\phi(x_i, \Theta) - c\|^2,$$

where  $c \in \mathcal{Z}$  is the one-class center,  $\Theta$  is the set of parameters of a representation network  $\phi$ . Deep SVDD obtains the sphere of the smallest volume by minimizing the  $\mathcal{L}_{svdd}$  in the representation space  $\mathcal{Z} \subset \mathbb{R}^K$ .

**Negative-sample-free contrastive learning.** BYOL [12], SimSiam [7], and Vicreg [2] are representatives of contrastive learning without negative pairs. The optimization objective of SimSiam is simplified as:

$$(3.2) \quad \mathcal{L}_{sim} = \frac{1}{N} \sum_{i=1}^N -\frac{z_i}{\|z_i\|_2} \cdot \frac{z'_i}{\|z'_i\|_2},$$

where  $z_i$  and  $z'_i$  are the representations of contrasting objects (positive pairs) in the latent space  $\mathcal{Z}$ . Equation (3.2) is essentially pulling the positive pair close using cosine similarity. As for the “hypersphere collapse”

caused by no negative pairs, BYOL and SimSiam solve it by bootstrap and asymmetric networks, and Vicreg by variance.

**Invariance term of COCA objective.** A crude way to integrate one-class classification and contrastive learning is treating it as multi-task learning with two adjustable hyper-parameters  $\alpha$  and  $\beta$  as follows:

$$(3.3) \quad \alpha \cdot \mathcal{L}_{svdd} + \beta \cdot \mathcal{L}_{sim}.$$

Therefore, the main intuition behind our model is that a positive correlation exists between one-class classification and contrastive learning, so their objectives can be achieved simultaneously by a loss function without hyper-parameters  $\alpha$  and  $\beta$ . Considering  $\text{sim}(u, v) = u^T v / \|u\|_2 \|v\|_2$  denotes cosine similarity between  $u$  and  $v$ , we define the invariance term  $d$  between  $\ell_2$ -normalized  $Q = \{q_1, q_2, \dots, q_N\}$  and  $Q' = \{q'_1, q'_2, \dots, q'_N\}$  as:

$$(3.4) \quad d(Q, Q') = \frac{1}{N} \sum_{i=1}^N \{[1 - \text{sim}(q_i, Ce)] + [1 - \text{sim}(q'_i, Ce)]\},$$

where  $Ce$  is the  $\ell_2$ -normalized one-class center defined by:

$$(3.5) \quad Ce(Q, Q') = \frac{1}{2N} \sum_{i=1}^N (q_i + q'_i).$$

Here,  $Ce$ ,  $q_i$ , and  $q'_i$  are distributed on the unit hypersphere after normalization. According to Equation (3.1)  $\mathcal{L}_{svdd}$ , minimizing  $d(Q, Q')$  brings  $q_i$  and  $q'_i$  closer to  $Ce$ , which achieves the one-class classification-based normality assumption. Meanwhile, on the unit hypersphere,  $d(Q, Q')$  and  $\mathcal{L}_{sim}$  are related as follows:

$$(3.6) \quad d(Q, Q') \geq 1 + \mathcal{L}_{sim}(Q, Q'),$$

which becomes tighter as  $d(Q, Q')$  decreases. Also, observe that minimizing the  $d(Q, Q')$  shrinks an upper bound of contrastive errors  $\mathcal{L}_{sim}(Q, Q')$ , and achieves the contrastive learning-based normality assumption. For more details see sub-section 3.4.

For the case where a little bit of training data is anomalous, which is very common in AD tasks, the *soft-boundary invariance* of the COCA objective employing the hinge loss function is defined as:

$$(3.7) \quad d_{soft}(Q, Q') = L + \frac{1}{vN} \sum_{i=1}^N \max\{0, S_i - L\},$$

where  $L$  is the  $(1-v)$ -quantile of  $S = \{S_1, S_2, \dots, S_N\}$ , hyper-parameter  $v \in (0, 1]$  controls the trade-off between  $L$  and violations of the boundary, i.e. the amount

of time series allowed to be mapped outside the boundary.  $S_i$  is the *anomaly score* of a time series  $\mathbf{X}_i$ , which is defined as:

$$(3.8) \quad S_i(\mathbf{X}_i) = 2 - \text{sim}(q_i, Ce) - \text{sim}(q'_i, Ce),$$

where,  $0 < S_i(\mathbf{X}_i) \leq 2$ .

**Variance term of COCA objective.** In AD, COCA removes negative pairs to avoid performance degradation caused by pushing away negative pairs that are both normal. However, both negative-sample-free contrastive learning and Deep SVDD are likely to give an undesired trivial solution that all outputs “collapse” to a constant, i.e. “hypersphere collapse”. Inspired by [2,8], COCA can then define the variance  $v$  as a hinge function on the standard deviation of the projected vectors  $q_i$ :

$$(3.9) \quad v(Q) = \frac{1}{N} \sum_{i=1}^N \max \left\{ 0, \gamma - \sqrt{\text{Var}(q_i)} + \varepsilon \right\},$$

where  $\gamma$  is a constant target value of the standard deviation, and  $\varepsilon$  is a small scalar to prevent instabilities. In our experiments,  $\gamma$  is set to 1, and  $\varepsilon$  is set to  $10^{-4}$ . On the other hand, according to the research in Deep SVDD, selecting an appropriate one-class center can alleviate the problem of hypersphere collapse. In COCA, the one-class center  $Ce$  is ensured to be non-zero in any dimension and only updated in the first few epochs, because experiments show that an unfixed  $Ce$  would make the network easily converge to a trivial solution.

The overall loss function of COCA is a weighted average of the invariance and variance terms:

$$(3.10) \quad \mathcal{L} = \lambda d(Q, Q') + \frac{\mu}{2}(v(Q) + v(Q')),$$

where  $\lambda$  and  $\mu$  are hyper-parameters controlling the contribution of each term in the loss. So similarly, the *soft-boundary* loss function of COCA is defined as:

$$(3.11) \quad \mathcal{L}_{soft} = \lambda d_{soft}(Q, Q') + \frac{\mu}{2}(v(Q) + v(Q')).$$

$\mathcal{L}$  applies to the training set without anomalies, while  $\mathcal{L}_{soft}$  is for those containing a few anomalies. Contrastive learning has two key properties: alignment and uniformity [28] (detail in 2). There is an inverse relationship between uniformity and hypersphere collapse, the better the uniformity the less likely the collapse will occur, and vice versa. Nevertheless, uniformity somewhat contradicts the aim of one-class classification [27], because the latter is to bring representations closer to the center on the unit hypersphere, while some representations may be instead pulled far away by uniformity.

Therefore, in our experiments,  $\lambda$  is fixed to 1, and  $\mu$  is determined by a grid search with the base condition  $\mu < 1$ .

**Anomaly Detection.** In the test phase, an anomaly score  $S_i$  will be generated for the time series  $\mathbf{X}_i$ . Then, the following formula is applied to determine whether  $\mathbf{X}_i$  can be classified as an anomaly:

$$(3.12) \quad x_t = \begin{cases} \text{anomaly}, & S_i > \tau \\ \text{normal}, & S_i \leq \tau \end{cases},$$

where  $\tau$  is a predefined threshold. The overall algorithm is summarized in supplementary material A.2.

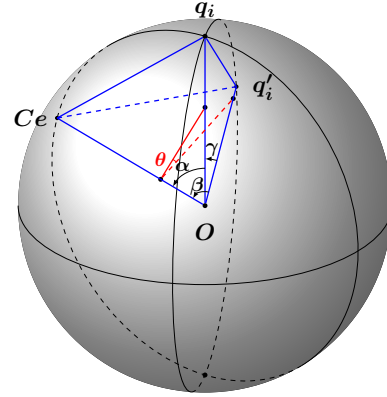


Figure 2: Invariance term schematic.  $O$  is the center of the unit hypersphere,  $Ce$  is the  $\ell_2$ -normalized one-class center,  $q_i$  and  $q'_i$  are  $\ell_2$ -normalized projected vectors,  $\theta$  is the dihedral angle between plane  $CeOq_i$  and  $CeOq'_i$ ,  $\alpha$  and  $\beta$  are one-class errors, and  $\gamma$  is the contrastive error.

**3.4 Relation to Contrastive Learning.** COCA treats representations  $q_i$  and reconstructed representations  $q'_i$  as positive pairs to learn shared information between different time steps of time series, discarding low-level information that is computationally expensive and unnecessary. Along with CPC [20], SimCLR [6], and wav2vec [1], though different in the types of positive pairs, COCA is essentially computing loss in the representation space. Therefore, maximizing the cosine similarity of  $q_i$  and  $q'_i$  in COCA is a type of negative-sample-free contrastive learning, and we name it “sequence contrast”. For time-series AD, COCA outperforms SimCLR-similar contrast methods that regard various augmentations as positive pairs (see sub-section 4.3).

Next, we will explain the mechanism of the invariance term to achieve the contrastive learning-based normality assumption. As shown in Fig. 2, on the unit

hypersphere, the angle  $\alpha/\beta/\gamma$  is proportional to the Euclidean distance  $l_{q_i Ce}/l_{q'_i Ce}/l_{q_i q'_i}$  between two points. According to the triangle inequality, the relationship between the three Euclidean distances is  $l_{q_i Ce} + l_{q'_i Ce} \geq l_{q_i q'_i}$ . Therefore we are minimizing the cosine similarity between  $q_i$ ,  $q'_i$ , and  $Ce$  in Equation (3.4), which is an upper bound on the sequence contrastive learning errors between  $q_i$  and  $q'_i$ .  $\alpha$ ,  $\beta$  and  $\gamma$  are related as follows:

$$(3.13) \quad \cos\gamma = \cos\alpha\cos\beta + \sin\alpha\sin\beta\cos\theta,$$

where  $\theta$  is the dihedral angle. According to Equation (3.13), when  $\alpha \rightarrow 0$  and  $\beta \rightarrow 0$ ,  $\cos\gamma \rightarrow 1$ . Therefore, Equation (3.6) becomes tighter as  $d(Q, Q')$  becomes smaller, which was also verified in our experiments. For more formal proof details see supplementary material A.1.

## 4 Experiments

This section presents the experimental setup, baselines, COCA variants, main results, and hyper-parameter analysis. The code is available at <https://github.com/ruiking04/COCA>.

**4.1 Experimental Setup. Datasets.** Given the findings in [29], this paper abandons the flawed time-series AD datasets, such as Yahoo, Numenta, and NASA, and employs AIOps and UCR to evaluate the proposed model. The datasets considered are as follows.

- AIOps challenge (**AIOps**)<sup>1</sup>. This consists of well-maintained business cloud KPIs from some large Internet companies and contains 29 univariate time-series sub-datasets.
- UCR time series anomaly detection (**UCR**)<sup>2</sup> [9]. This contains 250 univariate time-series sub-datasets from various fields.

Table 1 summarizes these datasets. The time series are partitioned into length- $T$  sequences by a sliding window with time-step  $Ts \leq T$ . The two datasets both have a large number of samples and few anomalies, which is a challenge for some AD models. This table shows the number of sequences in the training/validation/testing set, and the percentage of anomalous samples in the training/testing set. The training sets of UCR don't contain anomalies, so the models are trained using Equation (3.10), while *soft-boundary* loss function Equation (3.11) is used for AIOps.

Table 1: Summary of time-series anomaly detection datasets

	AIOps	UCR
Number of sub-datasets	29	250
Variables	1	1
Domain	Cloud KPIs	Various
Length $T$	16	64
Time step $Ts$	2	4
Total samples	2961039	4830858
Training/validation/testing	40%/10%/50%	24%/6%/70%
Training/testing anomaly	2.98%/1.92%	0%/0.71%

**Evaluation Metrics.** In most cases, time-series anomalies occur as continuous-time intervals rather than isolated points, leading to difficulty in quantifying the predicted anomaly label sequence. In recent years, many evaluation metrics for time-series AD have been proposed, such as NAB Score, Point-Adjusted (PA), Revised Point-Adjusted (RPA) metrics, etc., but these metrics may overestimate the performance of the AD algorithm [16]. To achieve a rigorous evaluation of time-series AD, this paper uses two metrics: accuracy metric [18] and affiliation metrics [15]. The *accuracy* =  $n/250$  is a metric specifically for the UCR dataset, where  $n$  is the number of correctly predicted sub-datasets. Each sub-dataset in UCR contains only one anomaly segment, so as long as the predicted anomaly is within the correct region, this sub-dataset is considered correctly predicted. Affiliation metrics calculate precision/recall/F1-score metrics based on the concept of “affiliation” between the ground truth and the prediction sets. Note that affiliation metrics on the entire dataset are weighted averages of affiliation metrics for each sub-dataset:

$$F1_{\text{entire}} = \sum_{i=1}^M \frac{k_i}{K} F1_i,$$

where  $M$  is the number of sub-datasets,  $K$  is the total number of anomaly segments for the entire dataset, and  $k_i$  is the number of anomaly segments of the  $i$ -th sub-dataset.

**4.2 Baselines and COCA Variants.** The proposed approach is compared against the following unsupervised and self-supervised anomaly detection methods.

*Traditional Anomaly Detection Baselines.* Three commonly used traditional anomaly detection baselines are adopted: One-class SVM (OC-SVM) [26], Isolation Forest (IF) [17], and Random Cut Forest (RCF) [14].

*Deep Anomaly Detection Baselines.* Then, four deep anomaly detection methods: Deep one-class

<sup>1</sup><https://github.com/NetManAIOps/KPI-Anomaly-Detection>

<sup>2</sup>[https://www.cs.ucr.edu/~eamonn/time\\_series\\_data\\_2018/](https://www.cs.ucr.edu/~eamonn/time_series_data_2018/)

(Deep SVDD) [24], Spectral Residual CNN (SR-CNN) [23], Deep Autoencoding Gaussian Mixture Model (DAGMM) [30], and LSTM Encoder-decoder (LSTM-ED) [19].

*Contrastive Learning Anomaly Detection Baselines.*

Finally, two contrastive learning baselines are set: Contrastive Predictive Coding Anomaly Detection (CPC-AD) [10, 20] and Time Series Temporal and Contextual Contrasting Anomaly Detection (TS-TCC-AD) [11, 27].

For Deep SVDD, we use Conv1D and LSTM to implement its autoencoder architecture to process time-series data. Although DAGMM is initially designed for tabular data, in [3] it is used for time-series data. CPC is originally a method for sequential data, treating images as a sequence of pixels, so the network structure does not need to be changed significantly when processing time-series data. For TS-TCC-AD based on [27], TS-TCC [11] is used to learn the representation of time series in the pre-training phase, and Deep SVDD is used for AD in the fine-tuning phase.

**COCA Variants.** Moreover, we include the following five COCA variants as baselines to demonstrate the effectiveness of individual components in COCA. *NoAug* removes the time-series augmentations of COCA. *NoOC* removes the one-class classification of COCA to optimize the similarity of representations  $q_i$  and reconstructed representations  $q'_i$ . *NoCL* removes the contrastive learning of COCA to optimize the similarity of representations and one-class center. The difference between the variant NoCL and Deep SVDD is that the former contains a learnable nonlinear projector  $p_\theta$  network and no pre-training. *NoVar* removes the variance term of COCA to optimize the similarity of representations and one-class center. *COCA-vi* treats different augmentations (jittering and scaling) as positive pairs for contrast learning, similar to SimCLR [6].

**Implementation Details.** The network structure of our proposed COCA consists of two parts: encoder and Seq2Seq. The encoder comprises 2-block temporal convolutional modules that each are followed by batch normalization, ReLU activation, and  $2 \times 2$  max-pooling. For the Seq2Seq, two identical three-layer LSTMs are employed with the same dropout rate at 0.45 as 1D-CNNs. As for optimizer, an Adam optimizer with a learning rate from  $1e-4$  to  $5e-4$ , weight decay of  $5e-4$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.99$  is adopted. On the AIOps dataset, after calculating the anomaly scores, COCA searches on anomaly sample rate  $p$  from 0.01% to 0.30% with step 0.01% to determine the optimal anomaly threshold  $\tau$ . The UCR sub-datasets each have only one anomaly segment, so COCA directly takes the largest anomaly score as an anomaly. In addition, for UCR we use the early stopping strategy, as the sub-datasets from

Table 2: Average affiliation F1-score(%) and accuracy (%) with standard deviation for anomaly detection on time-series datasets. The best results are in bold.

Datasets	AIOps	UCR	
Metric	Affiliation F1	Affiliation F1	Accuracy
<b>OC-SVM</b>	25.36	60.26	8.80
<b>IF</b>	33.24	59.40	37.60
<b>RCF</b>	34.48±0.30	58.36±0.59	38.67±0.68
<b>Deep SVDD</b>	38.23±0.65	37.19±1.35	7.60±1.73
<b>SR-CNN</b>	31.54±1.03	51.72±0.83	30.40±0.91
<b>DAGMM</b>	36.15±0.95	66.93±0.47	6.13±0.50
<b>LSTM-ED</b>	34.12±0.54	66.87±1.07	51.02±2.05
<b>CPC-AD</b>	35.36±1.87	48.65±1.92	6.37±0.53
<b>TS-TCC-AD</b>	31.91±2.05	44.27±1.38	0.56±0.27
<b>COCA</b>	<b>66.78±2.91</b>	<b>79.16±1.27</b>	<b>66.12±2.62</b>
<b>NoAug</b>	65.74±4.61	57.24±1.35	26.64±2.35
<b>NoOC</b>	51.49±5.96	62.33±2.05	33.96±2.61
<b>NoCL</b>	63.80±3.29	77.80±1.82	63.84±3.65
<b>NoVar</b>	65.90±2.45	78.82±1.60	65.16±2.81
<b>COCA-vi</b>	65.86±3.07	75.48±1.20	60.36±2.25

different domains vary in epochs to convergence. Each method is run 10 times to obtain the mean and standard deviation. Lastly, all the models are built with PyTorch 1.7 and Merlion 1.1.1<sup>3</sup> [3], and trained on an NVIDIA Tesla V100 GPU. See more details about augmentation and hyper-parameters in supplementary material B.2.

**4.3 Main Results.** We report affiliation F1-score and accuracy in Table 2. From the vertical view of the table, some methods perform poorly on AIOps because there are anomalous samples in the training set, which leads to high false negative rates. On the other hand, for the UCR dataset, methods such as OC-SVM, Deep SVDD, and TS-TCC-AD have higher F1-score but lower accuracy. That's because the accuracy metric is binary (anomaly found or not), and it indicates these methods' results are close to the correct range of ground-truth anomaly but do not fall within it.

From the horizontal view of the table, four conclusions can be drawn. First, in shallow methods, RCF with F1-score over 34% performs well and even outperforms some deep methods, showing that RCF methods are good baselines in time-series AD. Second, TS-TCC-AD and Deep SVDD each have a performance gap of over 7% on the two datasets, indicating that regardless of pre-training methods, the pre-training process itself limits the performance. It also further confirms that the pre-trained deep model limits the performance of two-staged AD methods. Third, DAGMM and LSTM-

<sup>3</sup><https://github.com/salesforce/merlion>



ED perform better than other deep baselines, indicating the normality assumptions of clustering and reconstruction are more relevant to the nature of AD. Last, the proposed COCA outperforms all baselines on both datasets, demonstrating the effectiveness and robustness of ensemble multiple normality assumptions.

Also, Table 2 shows the effectiveness of each component in our proposed COCA model. To be more specific, by analyzing the AD performance of the NoAug, augmentations improve the performance of AD on the two datasets, especially on UCR. The results of COCA, NoOC, and NoCL show that the combination of multiple normality assumptions can improve the performance of AD effectively. Meanwhile, the NoVar performs poorly compared to COCA, which makes it clear that the variance term of the COCA objective is important. The COCA-vi is 2% averagely lower than the COCA on the two datasets because it treats different augmentations as positive pairs and ignores temporal dependencies. Overall, the results of COCA are better than the five variants, indicating the effectiveness and necessity of each component in our model.

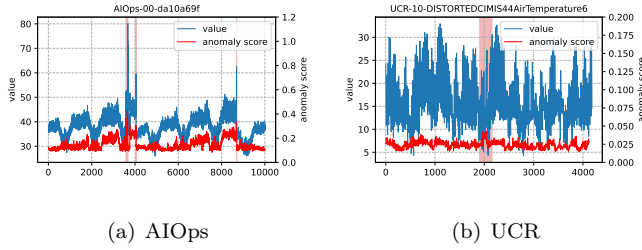


Figure 3: AD results of COCA on AIOps and UCR datasets. The blue curves are the original data value. The pink areas represent the ground-truth anomalies including point and subsequence anomalies. The red curves are the anomaly scores predicted by COCA.

**4.4 Visualization.** To provide a more intuitive evaluation, visualizations of AD on AIOps and UCR are conducted, in Fig. 3. It can be seen that AIOps contains many point anomalies, which are suitable for some AD methods that are specialized in learning global features. In contrast, UCR contains both point and sequence anomalies. COCA performs better on UCR compared to on AIOps, further illustrating that AD methods combining multiple normality assumptions can be applied to complex anomalous situations.

**4.5 Hyper-parameters Analysis.** In this section, sensitivity analysis is performed on the AIOps and UCR to study two main parameters:  $v \in (0, 1]$  in Equation (3.7) and the epoch  $e$  before stopping updating the

center  $C_e$ . Fig. 4(a) shows the effect of  $v$  on the overall performance, where the y-axis is the affiliation F1-score metric. For AIOps, we observe that  $v = 0.001$  is the best. Apparently, appropriate anomaly proportions  $v$  should be selected according to the anomaly proportion of datasets. Fig. 4(b) shows the results of varying epoch  $e$  of stopping update center  $C_e$  in a range between 1 and 50. The model is shown to perform best on UCR when  $e = 10$ , which suggests that the center  $C_e$  should be frozen early since updating the center  $C_e$  frequently increases the likelihood of hypersphere collapse.

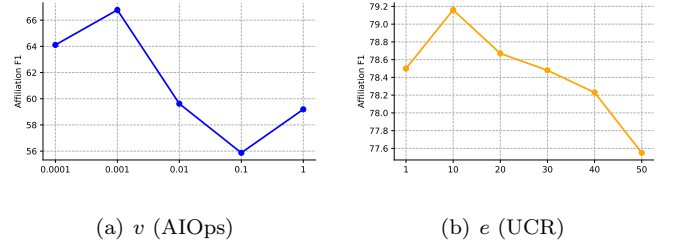


Figure 4: Two sensitivity analysis experiments on AIOps and UCR datasets. The left is the hyper-parameter  $v \in (0, 1]$  of *soft-boundary invariance* and the right is training epoch  $e$  before stopping updating the center  $C_e$ .

## 5 Conclusion

We propose a novel deep framework called COCA for unsupervised time-series anomaly detection. It combines the normality assumptions of contrastive learning and one-class classification, clarifies the essence of contrastive learning, and presents a new negative-sample-free type named “sequence contrast”. Specially, we present a novel contrastive one-class loss function optimizing the loss of both assumptions simultaneously in one stage without tuning hyper-parameters as in most multi-task learning, as well as preventing “hypersphere collapse”. Experiments on various datasets demonstrate that the performance of COCA achieves state-of-the-art. We hope our work can help deepen the understanding of contrastive learning and offer more possibilities for fusion studies of various anomaly detection methods.

## Acknowledgment

This research is supported by the Zhejiang Province Key R&D Program of China (2023C01070) and the Zhejiang Provincial Natural Science Foundation of China under Grant (LGG22F020043). The authors acknowledge the providers of datasets, including AIOps and UCR [9].

## References



- [1] A. BAEVSKI, Y. ZHOU, A. MOHAMED, AND M. AULI, *wav2vec 2.0: A framework for self-supervised learning of speech representations*, Advances in Neural Information Processing Systems, 33 (2020), pp. 12449–12460.
- [2] A. BARDES, J. PONCE, AND Y. LECUN, *Vicreg: Variance-invariance-covariance regularization for self-supervised learning*, in ICLR, 2022.
- [3] A. BHATNAGAR, P. KASSIANIK, C. LIU, T. LAN, W. YANG, R. CASSIUS, D. SAHOO, D. ARPIT, S. SUBRAMANIAN, G. WOO, ET AL., *Merlion: A machine learning library for time series*, arXiv preprint arXiv:2109.09265, (2021).
- [4] A. BLÁZQUEZ-GARCÍA, A. CONDE, U. MORI, AND J. A. LOZANO, *A review on outlier/anomaly detection in time series data*, ACM Computing Surveys (CSUR), 54 (2021), pp. 1–33.
- [5] M. CARON, I. MISRA, J. MAIRAL, P. GOYAL, P. BOJANOWSKI, AND A. JOULIN, *Unsupervised learning of visual features by contrasting cluster assignments*, in NeurIPS, 2020.
- [6] T. CHEN, S. KORNBLITH, M. NOROUZI, AND G. HINTON, *A simple framework for contrastive learning of visual representations*, in ICML, PMLR, 2020, pp. 1597–1607.
- [7] X. CHEN AND K. HE, *Exploring simple siamese representation learning*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 15750–15758.
- [8] P. CHONG, L. RUFF, M. KLOFT, AND A. BINDER, *Simple and effective prevention of mode collapse in deep one-class classification*, in IJCNN, IEEE, 2020, pp. 1–9.
- [9] H. A. DAU, E. KEOGH, K. KAMGAR, C.-C. M. YEH, Y. ZHU, S. GHARGHABI, C. A. RATANAMAHATANA, YANPING, B. HU, N. BEGUM, A. BAGNALL, A. MUEEN, AND H.-M. BATISTA, GUSTAVO, *The ucr time series classification archive*, October 2018.
- [10] P. DE HAAN AND S. LÖWE, *Contrastive predictive coding for anomaly detection*, arXiv preprint arXiv:2107.07820, (2021).
- [11] E. ELDELE, M. RAGAB, Z. CHEN, M. WU, C. K. KWONG, X. LI, AND C. GUAN, *Time-series representation learning via temporal and contextual contrasting*, IJCAI, (2021).
- [12] J.-B. GRILL, F. STRUB, F. ALTCHÉ, C. TALLEC, P. RICHEMOND, E. BUCHATSKAYA, C. DOERSCH, B. AVILA PIRES, Z. GUO, M. GHESLAGHI AZAR, ET AL., *Bootstrap your own latent-a new approach to self-supervised learning*, Advances in Neural Information Processing Systems, 33 (2020), pp. 21271–21284.
- [13] F. E. GRUBBS, *Procedures for detecting outlying observations in samples*, Technometrics, 11 (1969), pp. 1–21.
- [14] S. GUHA, N. MISHRA, G. ROY, AND O. SCHRIJVERS, *Robust random cut forest based anomaly detection on streams*, in ICML, PMLR, 2016, pp. 2712–2721.
- [15] A. HUET, J. M. NAVARRO, AND D. ROSSI, *Local evaluation of time series anomaly detection algorithms*, in ACM SIGKDD, 2022, pp. 635–645.
- [16] S. KIM, K. CHOI, H.-S. CHOI, B. LEE, AND S. YOON, *Towards a rigorous evaluation of time-series anomaly detection*, in AAAI, vol. 36, 2022, pp. 7194–7201.
- [17] F. T. LIU, K. M. TING, AND Z.-H. ZHOU, *Isolation forest*, in 2008 eighth IEEE international conference on data mining, IEEE, 2008, pp. 413–422.
- [18] Y. LU, R. WU, A. MUEEN, M. A. ZULUAGA, AND E. KEOGH, *Matrix profile xxiv: Scaling time series anomaly detection to trillions of datapoints and ultra-fast arriving data streams*, in ACM SIGKDD, 2022, pp. 1173–1182.
- [19] P. MALHOTRA, A. RAMAKRISHNAN, G. ANAND, L. VIG, P. AGARWAL, AND G. SHROFF, *Lstm-based encoder-decoder for multi-sensor anomaly detection*, arXiv preprint arXiv:1607.00148, (2016).
- [20] A. V. D. OORD, Y. LI, AND O. VINYALS, *Representation learning with contrastive predictive coding*, arXiv preprint arXiv:1807.03748, (2018).
- [21] G. PANG, C. SHEN, L. CAO, AND A. V. D. HENGEL, *Deep learning for anomaly detection: A review*, ACM Computing Surveys (CSUR), 54 (2021), pp. 1–38.
- [22] C. QIU, T. PFROMMER, M. KLOFT, S. MANDT, AND M. RUDOLPH, *Neural transformation learning for deep anomaly detection beyond images*, in ICML, PMLR, 2021, pp. 8703–8714.
- [23] H. REN, B. XU, Y. WANG, C. YI, C. HUANG, X. KOU, T. XING, M. YANG, J. TONG, AND Q. ZHANG, *Time-series anomaly detection service at microsoft*, in ACM SIGKDD, 2019, pp. 3009–3017.
- [24] L. RUFF, R. VANDERMEULEN, N. GOERNITZ, L. DEECKE, S. A. SIDDIQUI, A. BINDER, E. MÜLLER, AND M. KLOFT, *Deep one-class classification*, in ICML, PMLR, 2018, pp. 4393–4402.
- [25] T. SCHLEGL, P. SEEBÖCK, S. M. WALDSTEIN, U. SCHMIDT-ERFURTH, AND G. LANGS, *Unsupervised anomaly detection with generative adversarial networks to guide marker discovery*, in International conference on information processing in medical imaging, Springer, 2017, pp. 146–157.
- [26] B. SCHÖLKOPF, R. C. WILLIAMSON, A. J. SMOLA, J. SHAWE-TAYLOR, J. C. PLATT, ET AL., *Support vector method for novelty detection.*, in NIPS, vol. 12, Citeseer, 1999, pp. 582–588.
- [27] K. SOHN, C.-L. LI, J. YOON, M. JIN, AND T. PFISTER, *Learning and evaluating representations for deep one-class classification*, ICLR, (2021).
- [28] T. WANG AND P. ISOLA, *Understanding contrastive representation learning through alignment and uniformity on the hypersphere*, in ICML, PMLR, 2020, pp. 9929–9939.
- [29] R. WU AND E. KEOGH, *Current time series anomaly detection benchmarks are flawed and are creating the illusion of progress*, IEEE Transactions on Knowledge and Data Engineering, (2021).
- [30] B. ZONG, Q. SONG, M. R. MIN, W. CHENG, C. LUMEZANU, D. CHO, AND H. CHEN, *Deep autoencoding gaussian mixture model for unsupervised anomaly detection*, in ICLR, 2018.

## A Methodology Details

This section provides the details and hyper-parameters for COCA time-series AD.

**A.1 Estimating the Invariance Term.** As shown in Fig. 2, on the unit hypersphere, the formal proof is as follows:

$$\begin{aligned}
d(q_i, q'_i) &= [1 - \text{sim}(q_i, Ce)] + [1 - \text{sim}(q'_i, Ce)] \\
&\propto \alpha + \beta \\
&\propto l_{q_i Ce} + l_{q'_i Ce} \\
&= \sqrt{\|q_i - Ce\|^2} + \sqrt{\|q'_i - Ce\|^2} \\
&\geq \sqrt{\|q_i - q'_i\|^2} \\
&\propto \gamma \\
&\propto 1 - \text{sim}(q_i, q'_i) \\
&= 1 + \mathcal{L}_{sim}(Q, Q'),
\end{aligned}
\tag{A.1}$$

here,  $l_*$  are the Euclidean distances.  $\mathcal{L}_{sim}(Q, Q')$  is the contrastive error expressing the agreement between positive pairs.

**A.2 Detailed algorithms of COCA.** First, a pseudo-code for COCA in Pytorch style is provided in Algorithm 1.

**A.3 COCA Variants Loss Function.** Moreover, we include the following five COCA variants as baselines to demonstrate the effectiveness of individual components in COCA.

*NoAug.* The Variant NoAug removes the time-series augmentations of COCA.

*NoOC.* The Variant NoOC removes the one-class classification of COCA to optimize the similarity of representations  $q_i$  and reconstructed representations  $q'_i$ . Its invariance term of the loss function is defined as:

$$\frac{1}{N} \sum_{i=1}^N 1 - \text{sim}(q_i, q'_i).
\tag{A.2}$$

*NoCL.* The Variant NoCL removes the contrastive learning of COCA to optimize the similarity of representations and one-class center. Its invariance term of the loss function is defined as:

$$\frac{1}{N} \sum_{i=1}^N 1 - \text{sim}(q_i, Ce).
\tag{A.3}$$

The difference between the variant NoCL and Deep SVDD is that the former contains a learnable nonlinear projector  $p_\theta$  network and no pre-training.

---

## Algorithm 1 COCA's main training algorithm.

---

**Input:** a set of augmented time series (jittering and scaling)  $\{\mathbf{X}_i\}_{i=1}^N$ , batch size  $N$ , structure of  $f, g, h, p$ , constant  $nu, v, \gamma, \varepsilon, \lambda, \mu$ .

**Output:** Parameters of the network  $f, g, h$ , and  $p$ .

```

for sampled batch  $\{\mathbf{X}_i\}_{i=1}^N$  do
  for all  $i \in \{1, \dots, N\}$  do
    # representations
     $\mathbf{Z}_i = f(\mathbf{X}_i)$ 
     $q_i = p(\mathbf{Z}_i)$ 
    # reconstruction representations
     $\mathbf{Z}'_i = h(g(\mathbf{Z}_i))$ 
     $q'_i = p(\mathbf{Z}'_i)$ 
     $Ce = \frac{1}{2N} \sum_{i=1}^N (q_i + q'_i)$ 
    define  $\text{sim}(u, v)$  as  $\text{sim}(u, v) = u^T v / \|u\|_2 \|v\|_2$ 
    for all  $i \in \{1, \dots, N\}$  do
      # anomaly score
      define  $S_i(\mathbf{X}_i)$  as  $2 - \text{sim}(q_i, Ce) - \text{sim}(q'_i, Ce)$ 
    if soft-boundary then
       $L = \text{quantile}(S(\mathbf{X}), 1 - \eta)$ 
       $d(Q, Q') = L + \frac{1}{vN} \sum_{i=1}^N \max\{0, S_i - L\}$ 
    else
       $d(Q, Q') = \frac{1}{N} \sum_{i=1}^N S_i(\mathbf{X}_i)$ 
       $v(Q) = \frac{1}{N} \sum_{i=1}^N \max\{0, \gamma - \sqrt{\text{Var}(q_i) + \varepsilon}\}$ 
       $v(Q') = \frac{1}{N} \sum_{i=1}^N \max\{0, \gamma - \sqrt{\text{Var}(q'_i) + \varepsilon}\}$ 
       $\mathcal{L} = \lambda d(Q, Q') + \frac{\mu}{2} (v(Q) + v(Q'))$ 
      update networks  $f, g, h$ , and  $p$  to minimize  $\mathcal{L}$ 
    return network  $f, g, h$ , and  $p$ 

```

---

*NoVar.* The Variant NoVar removes the variance term of COCA to optimize the similarity of representations and one-class center. Its loss function is defined as:

$$d(Q, Q').
\tag{A.4}$$

*COCA-vi.* The variant COCA-vi treats different augmentations (jittering and scaling) as positive pairs for contrast learning, similar to SimCLR [6]. Its invariance term of the loss function is defined as:

$$d(Z^1, Z^2),
\tag{A.5}$$

where  $Z^1$  and  $Z^2$  are the representations of the time-series data after jittering and scaling, respectively.

## B Experiments.

**B.1 Baseline Details.** For these deep baselines, Table 3 shows the normality assumptions, study domains, and whether two-staged.

Table 3: Summary of deep baselines.

	Assumption	Two-staged	Original domain
Deep SVDD	Autoencoder&One-class	✓	Image
SR-CNN	Saliency map	×	Time series
DAGMM	Clustering	×	Tabular data
LSTM-ED	Autoencoder	×	Time series
CPC-AD	Contrast	×	Time series
TS-TCC-AD	Contrast&One-class	✓	Time series

**B.2 Hyper-parameters Details.** COCA is implemented in PyTorch, and some important parameter values used in the model are listed here, see Table 4. In this table, *repre\_channels* is the dimension of the final representations  $\mathbf{Z}$ , *hidden\_size* is the dimension of the Seq2Seq in the model, and *project\_channels* is the dimension of the projector. *window\_size* is the size of time window, the same as the length of time series  $T$ , and *time\_step* is the step while sliding. *stop\_change\_center* is the training epoch  $e$  before stopping updating the center  $Ce$ .  $\mu$  is the weight of the variance term of COCA objective.  $lr$  is the learning rate and  $nu$  is the hyper-parameter  $v \in (0, 1]$  of *soft-boundary invariance*. *scale\_ratio* and *jitter\_ratio* are the rate of scaling and jittering while applying data augmentation, respectively.

comparison error  $\text{sim}(q_i, q'_i)$ .

Table 4: The values of hyper-parameters used in COCA

	AIOps	UCR
<i>repre_channels</i>	32	64
<i>hidden_size</i>	64	128
<i>project_channels</i>	16	32
<i>window_size</i>	16	64
<i>time_step</i>	2	4
<i>stop_change_center</i>	1	10
$\mu$	0.1	0.1
<i>lr</i>	0.0001	0.0003
<i>nu</i>	0.001	-
<i>scale_ratio</i>	1.1	0.8
<i>jitter_ratio</i>	0.1	0.2

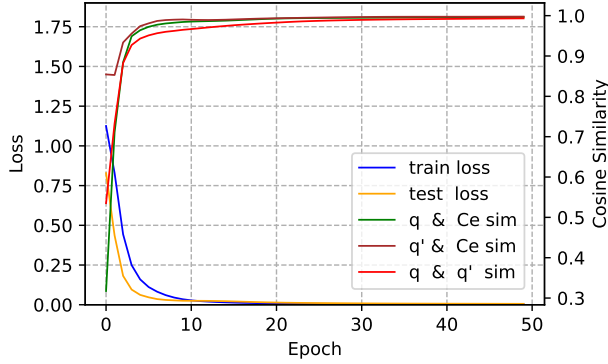


Figure 5: Loss and cosine similarity results for COCA and COCA-NoVar on UCR. Blue: train loss, Orange: test loss, Green:  $\text{sim}(q_i, Ce)$ , Brown:  $\text{sim}(q'_i, Ce)$ , Red:  $\text{sim}(q_i, q'_i)$ .

**B.3 Relation to Contrastive Learning.** To verify the validity of the invariance terms in the loss function of COCA, Fig. 5 illustrates loss and cosine similarity results for COCA on UCR. As can be seen from Fig. 5, the process of optimizing the loss function  $\mathcal{L}$  makes  $\text{sim}(q_i, Ce) \rightarrow 1$ ,  $\text{sim}(q'_i, Ce) \rightarrow 1$  and  $\text{sim}(q_i, q'_i) \rightarrow 1$ , which indicates that the loss we design not only makes  $q_i$  and  $q'_i$  closer to  $Ce$ , but also minimizes the sequence