# TRUSTED AND INTERACTIVE CLUSTERING FOR TIME-SERIES DATA

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Time-series clustering has gained abundant popularity and has been used in diverse scientific areas. However, few researchers take an information fusion perspective to combine information from the time and frequency domains to accomplish clustering, although these two domains offer distinct and complementary characteristics of time-series. Motivated by this issue, we propose a trusted and interactive model, which leverages evidence theory to combine time- and frequencybased clustering results produced by the corresponding contrastive learning module. After mathematizing clustering results from the two domains as mass functions, the uncertainty contained in these results can be quantified at the samplespecific level. The combined result thus promotes clustering reliability, and is optimized based on the pseudo-labels generated by k-means in an interactive learning paradigm. Both theoretical analysis and experimental results on 136 benchmark datasets validate the effectiveness of the proposed model in clustering performance. Extensive ablation experiments demonstrate the contribution of combining information from the time and frequency domains and using the interactive learning paradigm. The embeddings learned are also experimentally shown to perform well in other downstream tasks.

027 028 029

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

025

026

## 1 INTRODUCTION

Time-series clustering is an important data mining technology widely applied in different fields, such as sensor data analysis (Hayashi et al., 2024), anomaly detection (Middlehurst et al., 2024) and medical field (Zhang et al., 2024), aiming to segment time-series data samples into patterns (called clusters) with homologous characteristics (Gong et al., 2022). Unlike images, time-series data generally do not show human-recognizable features to different classes, because the label information is contained in not only the time domain but also the frequency domain (Zhang et al., 2022).

037 **Related work.** Most of the existing time-series clustering methods focus on the time domain, and 038 can be divided into two categories: raw-data-based methods and feature-based ones (Ma et al., 2022). Raw-data-based methods perform clustering based on a modified similarity metric, which 040 quantifies the distance more appropriately between time-series samples (Hayashi et al., 2024; Ferreira & Zhao, 2016; Paparrizos & Gravano, 2015; Yang & Leskovec, 2011). Feature-based methods 041 have recently garnered greater attention, because the raw-data-based ones are not capable of mod-042 eling nonlinear temporal dependencies and multiscale (long and short-term) temporal dependencies 043 (Ma et al., 2019a). Feature-based methods extract first informative features from raw samples in the 044 time domain, and then the clustering algorithms are conducted on the learned features (Tang et al., 045 2021; Fortuin et al., 2020). Some recent works (Zhang et al., 2024; Guijo-Rubio et al., 2021) also 046 optimize feature extraction and clustering jointly by introducing pseudo-label. For example, authors 047 in (Péalat et al., 2023) embed the time-series onto the Stiefel manifold to obtain the geometric rep-048 resentations of time-series samples. STCN (Ma et al., 2022) adopts a recurrent neural network and a self-supervised clustering module, which are trained iteratively by contributing to each other. Some recent works (e.g., (Fortuin et al., 2020; Tonekaboni et al., 2022)) leverage contrastive representation 051 learning for clustering time-series to further improve the performance. Few of the mentioned methods take an information fusion perspective to incorporate information from the time and frequency 052 domains to accomplish clustering, although these two domains offer distinct and complementary characteristics of time-series. More detailed related work is provided in Appendix A.1.

054 Motivation. In fact, the time domain depicts the temporal evolution of signal readouts, while the fre-055 quency domain reveals the distribution of signal magnitude across different frequency components 056 within the entire spectrum (Hyndman & Athanasopoulos, 2018). By explicitly incorporating the 057 frequency domain, a comprehension of time-series behavior can be attained, encompassing aspects 058 that are not fully captured by analyzing the time domain alone. Therefore, the first motivation for this work is incorporating frequency information to enhance the ability to detect clusters. Besides, the time and frequency domains can be regarded as distinct views of the same data (Cohen, 1995), 060 and they are interconvertible through Fourier and inverse Fourier Transformation (Brigham, 1988). 061 Given the temporal dynamics inherent in time-series data, the weights (i.e., quality) of the time and 062 frequency domains experience fluctuations over time. Then, the second motivation is to dynamically 063 describe the importance of clustering results derived from both the time and frequency domains. 064 Finally, to avoid poisoning the final result with the low-quality result from one domain, the last mo-065 tivation is to facilitate the appropriate integration of time- and frequency-based clustering results 066 under the fast-evolving uncertain scenario. 067

- Technical issues. To address the three challenges outlined in the motivation, three technical issues must be solved.
- (1) How to develop frequency-based contrastive augmentations to obtain clustering results?
- Despite the universal importance of frequency information in time-series and its pivotal role in classic signal processing (Soklaski et al., 2022), it is seldom explored in contrastive clustering for time-series data (Tonekaboni et al., 2022). This is because even a slight perturbation in the frequency domain could lead to significant alterations in the temporal patterns of the time domain (Flandrin, 1998).
- (2) How to quantify uncertainty in clustering results from the time and frequency domains?
  Quantifying the uncertainty is practical for enhancing the performance of such a "two-view" timeseries clustering, where the higher the uncertainty of a particular domain, the lower the weight
  contributing to the final result. Further, the uncertainty of every domain varies for different timeseries samples.
- (3) How to formulate and optimize the fusion of clustering results with uncertainty from the time and frequency domains? Fusion of the clustering results from the two domains belongs to the late fusion (Wang et al., 2019; Liu et al., 2018), most of which primarily cater to scenarios without any uncertainty (Wang et al., 2021). In our model, the fusion module assumes the critical responsibility of effectively managing uncertainty, all integrated seamlessly within an end-to-end framework.
- Contributions. We introduce a trusted clustering model (named TIC) for time-series data, integrat ing results from the time and frequency domains within an interactive learning paradigm (shown in
   Fig.1). In summary, the contributions of this paper are:
- We propose to adopt novel augmentations dedicated to clustering the frequency spectrum data, through contrastive sample discrimination. This could be the first work to leverage contrastive augmentation in the frequency domain in a time-series clustering problem.
- We represent the clustering results for each sample from the time and frequency domains as two distinct mass functions (Shafer, 1976), where the sample-specific uncertainty is accurately quantified within the framework of evidence theory (DEMPSTER, 1967). The Dirichlet distribution is used to mathematize the mass function and to model the class probability distribution from each domain.
- We propose an interactive framework for time-series clustering, which could be inspiring in multiple research areas of time-series. The trusted result (obtained by combining mass functions from the time and frequency domains) and pseudo-labels (derived from k-means) contribute to each other, allowing interactive model optimization.
- 102 103

090

## 2 PRELIMINARIES: EVIDENCE THEORY

104 105

106 Consider a variable  $\omega$  taking values in a finite set called the *frame of discernment*  $\Omega = \{\omega_1, \omega_2, \dots, \omega_C\}$ , where C is the number of clusters in a clustering problem. A mass function (also called a piece of evidence) (Shafer, 1976) is defined as a mapping from  $2^{\Omega}$  to [0,1] such that

122 123

124

125

141

142

143 144 145

146 147

 $\sum_{A \subseteq \Omega} m(A) = 1, m(A) > 0, \text{ where } 2^{\Omega} \text{ is the power set of } \Omega \text{ and } A \text{ denotes various subsets of } \Omega.$ These subsets are called the *focal sets* of *m*. The value of *m*(*A*) denotes a degree of belief assigned to the hypothesis " $\omega \in A$ ". The vacuous mass function verifies

$$m(\Omega) = m_{\Omega} = 1 \tag{1}$$

113 corresponding to total "ignorance", representing that  $\omega$  may belong to any subset of the  $\Omega$  (including the  $\Omega$  itself). In other words, the value of  $m(\Omega)$  measures the uncertainty about the values of 115 variable  $\omega$ . In this work, only the singleton focal sets  $\omega_1, \omega_2, \dots, \omega_C$  and ignorance focal set  $\Omega$ 116 are considered, i.e., the mass function is in the form of  $\mathbf{m} = [m(\omega_1), m(\omega_2), \dots, m(\omega_C), m(\Omega)]$ 117 (abbreviated as  $[m_1, m_2, \dots, m_C, m_\Omega]$ ).

118 Assume that there are 2 mass functions  $m_1$  and  $m_2$  on the same frame of discernment  $\Omega$ . The *Dempster's rule* (Shafer, 1976) used to pool the information provided by  $\mathbf{m}_1$  and  $\mathbf{m}_2$  is noted as  $\bigoplus$ and is defined by 121  $\sum_{n=1}^{\infty} m_n m_n$ 

$$m_{1\oplus 2}(A) = \frac{\sum_{B\cap C=A} m_{1B} m_{2C}}{\mathcal{K}_{12}}, \ A \neq \emptyset \& A \subseteq \Omega$$
<sup>(2)</sup>

where B and C are also the focal sets, and  $\mathcal{K}_{12} = 1 - \sum_{B \cap C = \emptyset} m_{1B} m_{2C}$  is a normalizing factor. Dempster's rule is commutative and associative (Shafer, 1976).



Figure 1: Overview of the proposed TIC model. TIC has three modules, a time-based contrastive module (colored yellow), a frequency-based contrastive module (colored blue), and an interactive learning module (colored red).

#### 3 THE PROPOSED METHOD: TIC

**Notions and problem formulation.** We are given a time-series dataset  $T = {\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n}$  of nunlabeled time-series samples, and sample  $\mathbf{x}_i$  has K channels and L time-steps. The goal is to group these n samples into C (a given value) clusters, which can be denoted by the *framework of discernment*  $\Omega = {\omega_1, \omega_2, \dots, \omega_C}$  (defined in Section 2). Without loss of generality, in the following, we focus on univariate (single-channel) time-series datasets, while noting that our TIC method can also accommodate multi-variate time-series. Superscript denotes contrastive augmentations,  $\mathbf{x}_i \equiv \mathbf{x}_i^{\text{Time}}$ denotes the input time-series and  $\mathbf{x}_i^{\text{Freq}}$  denotes the frequency spectrum of  $\mathbf{x}_i$ .

**Time-based contrastive module.** For the sample  $\mathbf{x}_i^{\text{Time}}$  in one mini-batch, a set of augmentations  $\mathcal{X}_i^{\text{Time}}$  is generated through a time-based augmentation bank  $\mathcal{B}^{\text{Time}} : \mathbf{x}_i^{\text{Time}} \to \mathcal{X}_i^{\text{Time}}$  including jittering, scaling, time-shifts and other common techniques (Eldele et al., 2021). Note that the augmentations in one mini-batch are produced using diverse techniques from the augmentation bank, to expose the model to complex temporal dynamics and obtain robust embeddings.

160 161  $\mathbf{x}_i^{\text{Time}}$  and a randomly selected augmentation  $\mathbf{\widetilde{x}}_i^{\text{Time}} \in \mathcal{X}_i^{\text{Time}}$  are fed into the encoder, denoted by  $H_{\mathrm{T}}(\cdot)$ . The corresponding embeddings  $\mathbf{g}_i^{\text{Time}} = H_{\mathrm{T}}(\mathbf{x}_i^{\text{Time}})$  and  $\mathbf{\widetilde{g}}_i^{\text{Time}} = H_{\mathrm{T}}(\mathbf{\widetilde{x}}_i^{\text{Time}})$  are obtained.

194

195 196 197

209

212 213

Intuitively, embedding  $\mathbf{g}_{i}^{\text{Time}}$  should be close to the embedding  $\mathbf{\widetilde{g}}_{i}^{\text{Time}}$ , but far away from the embeddings  $\mathbf{g}_{j}^{\text{Time}}$  and  $\mathbf{\widetilde{g}}_{j}^{\text{Time}}$  produced from another sample  $\mathbf{x}_{j}^{\text{Time}}$  in the same mini-batch. Therefore, the positive pair is  $(\mathbf{x}_{i}^{\text{Time}}, \mathbf{\widetilde{x}}_{i}^{\text{Time}})$  and the negative pairs are  $(\mathbf{x}_{i}^{\text{Time}}, \mathbf{x}_{j}^{\text{Time}})$  and  $(\mathbf{x}_{i}^{\text{Time}}, \mathbf{\widetilde{x}}_{j}^{\text{Time}})$ . The normalized temperature-scaled cross-entropy loss associated with  $\mathbf{x}_{i}^{\text{Time}}$  is defined as (Chen et al., 2020)

$$\mathcal{L}_{i}^{\text{Time}} = -\log \frac{\exp(\sin(\mathbf{g}_{i}^{\text{Time}}, \widetilde{\mathbf{g}}_{i}^{\text{Time}}))/\tau}{\sum_{\mathbf{x}_{i}} \mathbb{1}_{i \neq j} \exp(\sin(\mathbf{g}_{i}^{\text{Time}}, H_{\text{T}}(\mathbf{x}_{j}))/\tau)},\tag{3}$$

where  $\sin(\mathbf{a}, \mathbf{b}) = \mathbf{a}^{\mathrm{T}} \mathbf{b} / \|\mathbf{a}\| \|\mathbf{b}\|$  is the cosine similarity,  $\tau$  is a temporal hyperparameter to adjust scale,  $\mathbb{1}$  is an indicator function equaling to 1 when  $i \neq j$  and 0 otherwise, and  $\mathbf{x}_j$  denotes the sample (or its augmented sample) different from  $\mathbf{x}_i$  in the same mini-batch. Minimization of  $\mathcal{L}^{\mathrm{Time}}$  enforces encoder  $H_{\mathrm{T}}(\cdot)$  to bring embeddings w.r.t. positive pairs closer together, and push embeddings w.r.t. negative pairs farther apart.

**Frequency-based contrastive module.** Although the frequency spectrum is informative, few methods leverage the frequency-based contrastive augmentation for clustering time-series (Tonekaboni et al., 2022). In this module, we use the Fourier Transformation (Brigham, 1988) to generate the frequency spectrum  $\mathbf{x}_i^{\text{Freq}}$  for sample  $\mathbf{x}_i^{\text{Time}}$ .

As shown in (Flandrin, 1998; Zhang et al., 2022), a minor perturbation in the frequency domain can 180 lead to significant changes in the corresponding time domain. To mitigate this issue, we manipu-181 late the amplitude to generate frequency-based augmentation. More concretely, the augmentation 182 bank  $\mathcal{B}^{\text{Freq}}$ :  $\mathbf{x}_i^{\text{Freq}} \to \mathcal{X}_i^{\text{Freq}}$ , where  $\mathcal{X}_i^{\text{Freq}}$  is a set of frequency-based augmentations, includes the 183 upgrade or downgrade of amplitude. We randomly select  $\beta$  (the number of components to be manip-184 ulated) frequency components, and change each of their amplitudes from the original value Amporia 185 to  $\gamma \operatorname{Amp}_{orig}$ ,  $0 \leq \gamma < 1$  (downgrade) or to  $\gamma \operatorname{Amp}_{orig}$ ,  $\gamma > 1$ ,  $\operatorname{Amp}_{orig} < \gamma \operatorname{Amp}_{orig} \leq \operatorname{Amp}_{\max}$  (upgrade), where  $\gamma$  is a pre-defined coefficient and  $\operatorname{Amp}_{\max}$  is the maximum amplitude. Similar to the time-based contrastive module, the positive pair is  $(\mathbf{x}_i^{\operatorname{Freq}}, \widetilde{\mathbf{x}}_i^{\operatorname{Freq}})$  and the negative pairs are  $(\mathbf{x}_i^{\operatorname{Freq}}, \mathbf{x}_j^{\operatorname{Freq}})$  and  $(\mathbf{x}_i^{\operatorname{Freq}}, \widetilde{\mathbf{x}}_j^{\operatorname{Freq}})$ . 186 187 188 189

After generating an augmentation  $\widetilde{\mathbf{x}}_{i}^{\text{Freq}} \in \mathcal{X}_{i}^{\text{Freq}}$ , we feed the frequency spectrum  $\mathbf{x}_{i}^{\text{Freq}}$  and  $\widetilde{\mathbf{x}}_{i}^{\text{Freq}}$ into the encoder  $H_{\text{F}}(\cdot)$ , and obtain the embeddings  $\mathbf{g}_{i}^{\text{Freq}}$  and  $\widetilde{\mathbf{g}}_{i}^{\text{Freq}}$ . The frequency-based loss for sample  $\mathbf{x}_{i}^{\text{Freq}}$  is calculated as

$$\mathcal{L}_{i}^{\text{Freq}} = -\log \frac{\exp(\sin(\mathbf{g}_{i}^{\text{Freq}}, \widetilde{\mathbf{g}}_{i}^{\text{Freq}}))/\tau}{\sum_{\mathbf{x}_{j}} \mathbb{1}_{i \neq j} \exp(\sin(\mathbf{g}_{i}^{\text{Freq}}, H_{\text{F}}(\mathbf{x}_{j}))/\tau)}.$$
(4)

**Uncertainty quantification.** As shown in Fig.1, the embeddings  $\mathbf{g}_i^{\text{Time}}$  and  $\mathbf{g}_i^{\text{Freq}}$  are fed into the fully connected (FC) layers to obtain the clustering results in time and frequency domains. 199 Taking the time domain as an example, the  $FC_T$  converts the continuous embeddings to vectors 200  $\mathbf{p}_{\mathrm{T}} = [p_{\mathrm{T}}^{\mathrm{T}}, p_{\mathrm{T}}^{\mathrm{T}}, \cdots, p_{\mathrm{C}}^{\mathrm{T}}]$ , which describes the probability of C mutually exclusive events after being 201 normalized through the softmax operator, and can be regarded as the parameters of a multinomial 202 distribution (Bishop & Nasrabadi, 2006). By replacing these parameters with the parameters of a 203 Dirichlet distribution, the clustering result from the  $FC_T$  can be represented as a distribution over 204 possible *softmax* outputs instead of a point estimation of one *softmax* output (Sensoy et al., 2018). That is, a Dirichlet distribution parametrized over the  $\mathbf{p}_{T,i} = [p_{i1}^T, p_{i2}^T, \cdots, p_{iC}^T]$  represents the den-205 sity of such a probability assignment w.r.t to sample  $\mathbf{x}_i^{\text{Time}}$ . Therefore, it can model the second-206 207 order probabilities and uncertainty for the clustering result of  $\mathbf{x}_{i}^{\text{Time}}$  (Jsang, 2018). The definition of 208 Dirichlet distribution is given in Definition 1.

**Definition 1** The Dirichlet distribution is a probability density function for a categorical distribution *p*. It can be characterized by *C* parameters  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_C]$  and is given by:

$$Dir(\boldsymbol{p}|\boldsymbol{\alpha}) = \begin{cases} \frac{1}{B(\boldsymbol{\alpha})} \prod_{c=1}^{C} p_{c}^{\alpha_{c}-1} & \text{for } \boldsymbol{p} \in \mathcal{S}_{C}, \\ 0 & \text{otherwise,} \end{cases}$$
(5)

214 215 where  $S_C$  is the C-dimensional unit simplex  $S_C = \{\mathbf{p} | \sum_{c=1}^{C} p_c = 1, 0 \le p_1, p_2 \cdots, p_C \le 1\}$  and  $B(\alpha)$  is the C-dimensional multinomial beta function.

In the clustering problem investigated in this paper, Subjective logic (Jsang, 2018) is used to associate the parameters  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \cdots, \alpha_C]$  of a Dirichlet distribution with the output  $\mathbf{p} = [p_1, p_2, \cdots, p_C]$ , where the Dirichlet distribution is considered as the conjugate prior of the corresponding multinomial distribution (Bishop & Nasrabadi, 2006). Inspired from (Sensoy et al., 2018), the parameter  $\alpha_c$  is calculated as

$$\alpha_c = p_c + 1. \tag{6}$$

Then, the mass function  $\mathbf{m} = [m_1, m_2, \cdots, m_C, m_\Omega]$  is determined as

$$m_c = \frac{p_c}{S} = \frac{\alpha_c - 1}{S} \qquad c = 1, 2, \cdots, C,$$
  
$$m_\Omega = \frac{C}{S},$$
(7)

where  $S = \sum_{c=1}^{C} (p_c + 1) = \sum_{c=1}^{C} \alpha_c$  is the Dirichlet strength (Jsang, 2018), and the value of  $m_{\Omega}$  quantifies the uncertainty of the clustering result (as discussed in Section 2). From Eq.(7), it can be inferred that the higher value of  $p_c$ , the more belief mass is assigned to  $m_c$ . Besides, the lower value of the sum of  $p_c$ , the higher uncertainty  $m_{\Omega}$  for the clustering result. For clarity, we give a specific example to explain the above formulation.

**Combining the clustering results.** After computing the mass functions  $\mathbf{m}_i^{\text{Time}}$  and  $\mathbf{m}_i^{\text{Freq}}$  from time and frequency domains, we use Dempster's rule to combine these two mass functions

$$\mathbf{m}_{i}^{\text{Comb}} = \mathbf{m}_{i}^{\text{Time}} \oplus \mathbf{m}_{i}^{\text{Freq}},\tag{8}$$

where  $\mathbf{m}_i^{\text{Comb}}$  is the combined mass function w.r.t. sample  $\mathbf{x}_i$ . According to Eq.(2), the specific calculation is

$$m_{ic}^{\text{Comb}} = \frac{m_{ic}^{\text{Time}} \cdot m_{ic}^{\text{Freq}} + m_{i\Omega}^{\text{Time}} \cdot m_{ic}^{\text{Freq}} + m_{ic}^{\text{Time}} \cdot m_{i\Omega}^{\text{Freq}}}{1 - \sum_{r \neq v} m_{ir}^{\text{Time}} \cdot m_{iv}^{\text{Freq}}}$$

$$m_{i\Omega}^{\text{Comb}} = \frac{m_{i\Omega}^{\text{Time}} \cdot m_{i\Omega}^{\text{Freq}}}{m_{i\Omega}^{\text{Time}} \cdot m_{i\Omega}^{\text{Freq}}}$$
(9)

243 244 245

246

247 248

240 241 242

221 222

223 224 225

226 227 228

229 230

231

232

233

234

235 236 237

 $m_{i\Omega}^{\text{comb}} = \frac{m_{i\Omega}}{1 - \sum_{r \neq v} m_{ir}^{\text{Time}} \cdot m_{iv}^{\text{Freq}}}.$ 

According to Eq.(7), the combined output for  $\mathbf{x}_i$  is calculated as

$$p_{ic}^{\text{Comb}} = m_{ic}^{\text{Comb}} \cdot S_i^{\text{Comb}} \text{ and } \alpha_{ic}^{\text{Comb}} = p_{ic}^{\text{Comb}} + 1,$$
(10)

where  $S_i^{\text{Comb}} = \frac{C}{m_{i\Omega}^{\text{Comb}}}$ . The time- and frequency-based clustering results are appropriately fused using Subjective logic and Dempster's rule.

252 **Remark 1. More intuitions: why the combined clustering results are trusted?** The produced mass value  $m_{i\Omega}^{\text{Comb}}$  allows the TIC model to assess the reliability of clustering results so as to avoid risky decisions. Besides, Dempster's rule has the following advantages: (1) the mass function 253 254  $\mathbf{m}^{\text{Comb}}$  obtained by combining a certain mass function (with small  $m_{\Omega}$ ) with an uncertain one (with 255 big  $m_{\Omega}$ ) is still certain. This means that as long as the clustering results from either domain are 256 trusted, then the combined clustering results are trusted, even if the results from the other domain 257 have significant uncertainty. (2) the  $\mathbf{m}^{\text{Comb}}$  obtained by combining two uncertain mass function 258 (with large  $m_{\Omega}$ ) remains uncertain. This means that if the clustering results from both the time and 259 frequency domains are untrusted, the combined clustering results are necessarily untrusted. And 260 users can abandon the results to avoid risks when facing this case. In order to analyze theoretically 261 these advantages, we give the following mathematized propositions, where the combined mass func-262 tion, the mass functions from the time and frequency domains are abbreviated as  $\mathbf{m}^{\mathrm{Co}}$ ,  $\mathbf{m}^{\mathrm{T}}$  and  $\mathbf{m}^{\mathrm{F}}$ . 263 Since  $\mathbf{m}^{\mathrm{T}}$  and  $\mathbf{m}^{\mathrm{F}}$  are equally important in the combination, the following Propositions still hold despite the exchange of the superscripts  $^{T}$  and  $^{F}$ . The corresponding proofs are shown in Appendix 264 A.2. 265

**Proposition 1** A large  $m_{\Omega}^{\mathrm{T}}$  does not lead to a large  $m_{\Omega}^{\mathrm{Co}}$ , when one of  $m_{c}^{\mathrm{F}}$  is large and  $m_{\Omega}^{\mathrm{F}}$  is small. In particular,  $\mathbf{m}^{\mathrm{Co}}$  is identical to  $\mathbf{m}^{\mathrm{T}}$ , if  $\mathbf{m}^{\mathrm{F}}$  is totally uncertain (i.e.,  $m_{\Omega}^{\mathrm{F}} = 1$ ).

269

**Proposition 2** The  $m_{\Omega}^{\text{Co}}$  is monotonically increasing with  $m_{\Omega}^{\text{T}}$  and  $m_{\Omega}^{\text{F}}$ .

275

281

282

283 284

285

287

289

290

291

292

293

295 296 297

298

299

300 301

302

303

305

306

307

270 Interactive learning module. As shown in the lower left part of Fig.1, we concatenate the embed-271 dings  $[\mathbf{g}^{\text{Time}}; \mathbf{g}^{\text{Freq}}]$  produced from the encoders  $H_{\text{T}}(\cdot)$  and  $H_{\text{F}}(\cdot)$ . These concatenated embeddings 272 of n samples are fed into k-means every t epochs and update the pseudo-labels to calculate the 273 interactive loss  $\mathcal{L}^{\text{Inte}}$ .

For sample  $\mathbf{x}_i$ , its one-hot pseudo-label vector is denoted as  $\mathbf{y}_i$  with  $y_{ic} = 1$  and  $y_{iv} = 0$  for all  $v \neq c$ . In the interactive learning module, we modify the conventional cross-entropy  $\mathcal{L}_i^{ce} = -\sum_{c=1}^C y_{ic} \log(p_{ic}^{\text{Comb}})$  as

$$\mathcal{L}_{i}^{'ce} = \int \left[\sum_{c=1}^{C} -y_{ic} \log(p_{ic})\right] \frac{1}{B(\alpha_{i})} \prod_{c=1}^{C} p_{ic}^{\alpha_{ic}-1} d\mathbf{p}_{i} = \sum_{c=1}^{C} y_{ic} \left(\psi(S_{i}) - \psi(\alpha_{ic})\right)$$
(11)

where  $\psi(\cdot)$  is the digamma function,  $S_i$  is the Dirichlet strength w.r.t.  $\mathbf{x}_i$  and we omit superscript <sup>Comb</sup> for brevity. Such a modification enforces the parameters  $\alpha_i^{\text{Comb}}$  of the combined Dirichlet distribution to be optimized based on the pseudo-label vector  $\mathbf{y}_i$ , i.e., enforces a large  $p_{ic}^{\text{Comb}}$  to be produced from the  $y_{ic} = 1$  in  $\mathbf{y}_i$ . To further shrink the  $p_{iv}^{\text{Comb}}$  w.r.t.  $y_{iv}$  to 0, the following KL divergence is considered

$$KL[Dir(\mathbf{p}_{i}|\widetilde{\boldsymbol{\alpha}}_{i})\|Dir(\mathbf{p}_{i}|\mathbf{1})] = \log\left(\frac{\Gamma(\sum_{c=1}^{C}\widetilde{\alpha}_{ic})}{\Gamma(C)\prod_{c=1}^{C}\Gamma(\widetilde{\alpha}_{ic})}\right) + \sum_{c=1}^{C}(\widetilde{\alpha}_{ic}-1)[\psi(\widetilde{\alpha}_{ic})-\psi(\sum_{v=1}^{C}\widetilde{\alpha}_{iv})],$$
(12)

where  $\widetilde{\alpha}_i = \mathbf{y}_i + (1 - \mathbf{y}_i) \odot \boldsymbol{\alpha}_i$  is the adjusted parameters of Dirichlet distribution, 1 is quite a flat Dirichlet distribution and  $\Gamma(\cdot)$  is the gamma function. Thus, the interactive loss for  $\mathbf{x}_i$  is

> $\mathcal{L}_{i}^{\text{Inte}} = \mathcal{L}_{i}^{'ce} + \iota KL[Dir(\mathbf{p}_{i}|\widetilde{\boldsymbol{\alpha}}_{i}) \| Dir(\mathbf{p}_{i}|\mathbf{1})],$ (13)

where  $\iota$  is the balance hyperparameter. In summary, the sample-specific loss of TIC model is

$$\mathcal{L}_{i}^{\mathrm{TIC}} = \mathcal{L}_{i}^{\mathrm{Inte}} + \mathcal{L}_{i}^{\mathrm{Time}} + \mathcal{L}_{i}^{\mathrm{Freq}}.$$
 (14)



304 Figure 2: Illustration of the learning process in TIC. When minimizing contrastive loss  $\mathcal{L}^{\text{Time}}$  and  $\mathcal{L}^{\text{Freq}}$ , the learned embeddings ( $\mathbf{g}^{\text{Time}}$  and  $\mathbf{g}^{\text{Freq}}$ ) affect the output of the FC layers (trusted clustering result), as shown by the blue arrow. When minimizing  $\mathcal{L}^{Inte}$ , the TIC model enforces the trusted clustering result to be "close" to the pseudo-labels, which are generated by inputting the embeddings 308 into k-means. In the way of backward propagation, the trusted clustering result also affects the 309 embedding learning, as shown by the red arrow. The parameter learning in FC layers and encoders  $(H_{\rm T}(\cdot), H_{\rm F}(\cdot))$  contribute to each other interactively. 310

311 Remark 2. More intuitions about the interactive process. In Fig.2, we show the three main com-312 ponents in TIC model. The embedding learned affects the trusted clustering result when minimizing 313 the contrastive loss  $\mathcal{L}^{\text{Time}}$  and  $\mathcal{L}^{\text{Freq}}$ . The pseudo-labels are produced by inputting the embeddings 314 into k-means and are considered when minimizing  $\mathcal{L}^{\text{Inte}}$ . In this way, the trusted clustering result 315 affects the embedding learning when performing backward propagation. This allows the parameter 316 learning of FC layers and encoders to contribute to each other interactively. Besides, considering 317 the interactive loss avoids the class collision issue, where each sample is identified as a cluster in 318 the embedding space (shown in Fig.3(a)) (Arora et al., 2019). This is because every positive pair 319 consists only of the sample and its augmentation, without considering any other samples that may 320 belong to the same cluster, when calculating the contrastive losses. The pseudo-labels in interactive loss are generated by k-means, where the concatenated embeddings  $[\mathbf{g}^{\text{Time}}; \mathbf{g}^{\text{Freq}}]$  are treated as 321 the input. Therefore, some basic information about clusters is included in the interactive loss and 322 improves the clustering performance (shown in Fig.3(b)). We demonstrate quantitatively the above 323 conclusion in Fig.4(b).

#### 4 **EXPERIMENTS**

324

325 326

327 328

330

331

332

333

334 335

344

345

347 348

364

In this section, we conduct some experiments to answer the following research questions (RQs):

- RQ1 (Comparison experiments): How does the clustering performance of TIC compare with that of other state-of-the-art time-series clustering algorithms?
- **RQ2** (Ablation study): How do the various components of TIC contribute to its performance?
- RO3 (Embedding evaluation): How about using the learned embedding for other downstream tasks (e.g., classification, anomaly detection)?
- **RQ4** (Hyperparameter sensitivity): what about the hyperparameter sensitivity of TIC?



Figure 3: Illustration of the embedding space with (a) and without (b) the interactive loss. Because only the augmentation of each sample is considered when determining positive pairs, it may cause the class collision issue in the embedding space (shown in (a)), i.e., each cluster consists of only one sample (Arora et al., 2019). Optimizing the  $\mathcal{L}^{\text{Inte}}$  allows k-means to provide some basic information 346 about the clusters and TIC to achieve better performance (shown in (b)).

Benchmark datasets. We use 136 benchmark (widely used in related work e.g., (Paparrizos & Gra-349 vano, 2015; Zhang et al., 2022)) time-series datasets to evaluate the TIC model. Eight datasets are 350 collected from the real-world and the remaining 128 ones are from the UCR database (Chen et al., 351 2015). Eight real-world datasets are multi-variate or univariate, and their application scenarios in-352 clude EEG and ECG analyses, and mechanical faulty detection. The description of these benchmark 353 datasets is shown in Table 1. More details about the datasets are shown in Appendix A.3. 354

355 Table 1: Specific description for the 136 benchmark datasets. The detailed description of the 128 356 UCR datasets can be found in (Chen et al., 2015). #Sam, #Clu and #Cha denote the number of 357 samples, clusters and channels, respectively.

Dataset	#Sam	#Clu	#Cha	Length
EMG	204	3	1	1,500
ECG	43,673	4	1	1,500
HAR	10,299	6	9	128
Gesture	560	8	3	315
FD-A	8,184	3	1	5,120
FD-B	13,640	3	1	5,120
SleepEEG	371,055	5	1	200
Epilepsy	11,500	2	1	178
UCR (128 datasets)	[40, 16,637]	[2, 60]	1	[15, 2,844]

365 **Baselines.** To answer RQ1, we consider the following 10 time-series clustering methods, i.e., **DTCR** 366 (Ma et al., 2019b), k-shape (Paparrizos & Gravano, 2015), SOM-VAE (Fortuin et al., 2020), STCN 367 (Ma et al., 2022), TMEK (Tang et al., 2021), TNC (Tonekaboni et al., 2022), TS3C<sub>ch</sub> (Guijo-Rubio 368 et al., 2021), UMAP (Péalat et al., 2023), USSL (Zhang et al., 2019) and VLSC (Duan & Guo, 2023) 369 in the comparison experiment. 370

To answer RQ3, we evaluate the embeddings learned by TIC on the other downstream tasks (i.e., 371 classification and anomaly detection). The included baselines are 8 unsupervised representation 372 learning methods, i.e., activity2vec (Aggarwal et al., 2019), BTSF (Yang & Hong, 2022), CLOCS 373 (Kiyasseh et al., 2021), MHCCL (Meng et al., 2023), TFC (Zhang et al., 2022), Triplet (Franceschi 374 et al., 2019), TS2Vec (Yue et al., 2022) and TSTCC (Eldele et al., 2022). More details of the 375 baselines are shown in Appendix A.4. 376

Implementation details. In our TIC model, we use two 2-layer Transformer (Vaswani et al., 2017) 377 as backbones for encoders  $H_{\rm T}(\cdot)$  and  $H_{\rm F}(\cdot)$ . FC<sub>T</sub> and FC<sub>F</sub> contain three fully-connected layers with 378 hidden dimensions  $d_1 = L$  (time-series length),  $d_2 = 128$  and  $d_3 = C$  (number of clusters), where 379 the softmax layer is replaced with the RELU to ensure that the network outputs are non-negative 380 values. These two FC modules do not share any parameters. The full spectrum (symmetrical) is used to guarantee that  $\mathbf{x}^{\text{Time}}$  and  $\mathbf{x}^{\text{Freq}}$  have the same number of dimensions, when transforming the time 381 382 domain to frequency domain. The Adam optimizer with a learning rate of  $\{0.0001, 0.0002, 0.0003\}$ and 2-norm penalty coefficient of 0.0005 is used. We use the batch size of  $\{8, 16, 32, 64, 128\}$ according to the dataset size, and use the training epoch of 100. We set  $\beta = 1, \gamma \in \{0.5, 1.2\}$  for 384 the frequency augmentation and  $\tau = 0.2$  in loss functions (3) and (4). The balance hyperparameter 385  $\iota$  in Eq.(13) is gradually increased to prevent TIC model from paying too much attention to the KL 386 divergence in the initial training stage. The pseudo-labels are updated every t = 20 epochs. We 387 use the code provided in the corresponding paper, and the hyperparameters are finely tuned within 388 the configuration provided therein based on the unsupervised metric Davies-Bouldin Index for fair 389 comparison. The supervised ARI, NMI, ACC metrics are considered. All models are implemented 390 with PyTorch on an NVIDIA A100 Tensor Core GPU. In Appendix A.5, the statistical comparison 391 result derived from the Friedman test and Nemenyi test is provided.

**RQ1. Comparison experiment.** We show the comparison results between TIC and the 10 timeseries clustering baselines in Table 5, where the ARI clustering metric are used (Rand, 1971). We recall that the larger the metrics, the better the clustering performance. After being fed to the correct number C of clusters, each algorithm is run 5 times and the results are recorded in the form of mean<sub>±std.deviation</sub>. In particular, the average ARI and the corresponding average std. deviations on the 128 UCR datasets are reported. The results w.r.t. NMI, ACC and running time are provided in Appendix A.6.

Table 2: ARI of different algorithms on benchmark datasets. The  $\bullet/\circ$  indicates whether TIC is statistically superior/inferior to a certain comparing baseline based on the paired *t*-test at a 0.05 significance level. The statistics of win/tie/loss are shown in the last row of each sub-table. The best and the second-best results, between which the performance gaps are shown in the row named "gap" in each sub-table, are colored blue and red.

ARI	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
DTCR	.782±.02●	.812±.01•	.584±.02•	.874±.03●	<u>.882</u> ±.02●	.801±.01•	.571±.01●	.872±.02●	.634±.0
k-shape	$.684_{+.02}$ •	$.578_{+.01}$ •	$.765 \pm .03$	$.723 + .02 \bullet$	$.860 + .01 \bullet$	<u>.811</u> +.02•	$.806 + .01 \bullet$	$.909_{+.02}$ •	.802 + .0
SOM-VAE	.845±.03•	.812±.02•	.774±.01	<u>.921</u> ±.01	.824±.03•	.732±.02•	.816±.02•	.931±.01•	.813±.0
STCN	$\underline{1}_{\pm 0}$	.730±.02●	.669±.02•	.825±.01●	.879±.02●	.803±.01•	.815±.02●	.770±.01●	$.541 \pm .0$
TMEK	$.651_{+.02}$ •	$.706 + .02 \bullet$	$.690_{+.01}$ •	$.807 + .01 \bullet$	$.682 + .02 \bullet$	$.782_{+.01}$ •	$.741 + .02 \bullet$	$.901 + .02 \bullet$	$.642_{+.0}$
TNC	.751±.01•	.805±.02•	.693±.01•	.851±.02•	.780±.03•	.593±.01•	.741±.02•	$.725 \pm .01 \bullet$	$.707 \pm .0$
TS3C <sub>ch</sub>	$.703 + .02 \bullet$	<u>.813</u> +.02•	$.706_{+.01}$ •	$.852_{+.01}$ •	$.711 + .02 \bullet$	$.682_{+.01}$ •	<u>.865</u> +.01	$.939_{+.02}$ •	<u>.851</u> +.0
UMAP	.859±.01•	.791±.01●	$.643 \pm .02 \bullet$	.852±.02•	$.725 \pm .01 \bullet$	.785±.01•	$.863 \pm .01$	.972±.02	$.848 \pm .0$
USSL	<u>.876</u> ±.01●	.745±.01●	.621±.02•	.597±.02●	.622±.03•	.592±.03●	.802±.02•	.925±.02●	.710±.0
VLSC	1 + 0	$.805 + .01 \bullet$	$.611_{+.02}$ •	$.823 + .02 \bullet$	$.654 + .02 \bullet$	$.498_{+.01}$ •	$.796_{+.01}$ •	$.942 \pm .02$	$.757_{+.0}$
TIC (ours)	$\underline{1}_{\pm 0}$	.879±.01	.796±.01	<u>.931</u> ±.01	.939±.01	.853±.02	.896±.01	<u>.980</u> ±.01	.883±.0
gap	.124	.066	.022	.010	.057	.042	.031	.008	.032
win/tie/loss	8/2/0	10/0/0	8/2/0	9/1/0	10/0/0	10/0/0	8/2/0	8/2/0	9/1/0

412 413 414

Overall, our TIC model wins 80 and has tied performance on 10 out of 90 trials, when it is statis-415 tically compared with 10 baselines based on three metrics. On average, our TIC model claims a 416 large performance gap of 0.043 over the best baselines. Concretely, the largest performance gap is 417 0.0124 on the EMG dataset and the smallest one is 0.008 on the Epilepsy dataset. Only on the EMG 418 dataset, STCN and VLSC yield the totally correct clustering result, and achieve the equal ARI with 419 TIC. One potential explanation is that EMG is a simple dataset with only 3 clusters and 204 samples, 420 and thus it may be easy to be modeled. On the ECG dataset, TIC outperforms the strongest base-421 lines by a large margin of 0.066. Because ECG includes four clusters consisting of 43,673 samples 422 that lead to a more complex clustering task, a combination of clustering information in the time and 423 frequency domains results in better performance. In summary, the better performance of TIC model can be attributed to (1) Both time-domain and frequency-domain contrastive losses are considered; 424 (2) Quantification of sample-specific uncertainty reduces the wrong assignment in clustering deci-425 sion; (3) Dempster's rule appropriately combines the time- and frequency-based clustering results; 426 (4) the class collision issue can be improved by including the interactive loss  $\mathcal{L}^{\text{Inte}}$  in  $\mathcal{L}^{\text{TIC}}$ . The 427 ablation study (more results shown in Appendix A.7) experimentally demonstrates the four points 428 above. 429

430 **RQ2.** Ablation study. We conduct ablation studies to evaluate the importance of every com-431 ponent in the developed TIC model. Concretely, we compare TIC model with its 3 variants: W/o  $\mathcal{L}^{\text{Time}}/\mathcal{L}^{\text{Freq}}$ : the loss function  $\mathcal{L}^{\text{Time}}/\mathcal{L}^{\text{Freq}}$  is removed from the  $\mathcal{L}^{\text{TIC}}$ , i.e., the timebased/frequency-based contrastive module is removed; W/o  $\mathcal{L}^{\text{Inte}}$ : the loss function  $\mathcal{L}^{\text{Inte}}$  is removed from the  $\mathcal{L}^{\text{TIC}}$ , i.e., *k*-means does not provide pseudo-labels for the training.

Table 3: ARI values (mean $\pm$ std.deviation) of different variants of TIC.

ARI	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR	Average
W/o $\mathcal{L}^{\mathrm{Time}}$	$.958 \pm .01$	$.815 \pm .02$	$.702 \pm .01$	$.846 \pm .02$	$.859 \pm .02$	$.745 \pm .02$	$.687 \pm .02$	$.899 \pm .01$	$.756_{\pm.01}$	0.807
W/o $\mathcal{L}^{\mathrm{Freq}}$	$.934 \pm .01$	$.827 \pm .01$	$.754 \pm .01$	$.798 \pm .02$	$.547 \pm .02$	$.609 \pm .01$	$.781 \pm .02$	$.854 \pm .01$	$.716 \pm .01$	0.758
W/o $\mathcal{L}^{ ext{Inte}}$	$.042 \pm .03$	$.121 \pm .02$	$.098 \pm .01$	$.057 \pm .03$	$.069 \pm .02$	$.210_{\pm.01}$	$.009 \pm .02$	$.147 \pm .02$	$.059 \pm .01$	0.090
TIC (Full model)	$\underline{1}_{\pm 0}$	<u>.879</u> ±.01	<u>.796</u> ±.01	<u>.931</u> ±.01	<u>.939</u> ±.01	<u>.853</u> ±.02	<u>.896</u> ±.01	<u>.980</u> ±.01	<u>.883</u> ±.01	<u>0.906</u>

The results of the ablation study are reported in Table 3. As can be seen, removing  $\mathcal{L}^{\text{Time}}$  and  $\mathcal{L}^{\text{Freq}}$ 442 leads to performance degradation (average ARI) of 0.912 - 0.807 = 0.105 and 0.912 - 0.758 =443 0.154, respectively. In particular, on the datasets {FD-A, FD-B} with a high sampling frequency of 444 64k Hz, removing  $\mathcal{L}^{\text{Freq}}$  cause more pronounced performance degradation, e.g., the ARI values of 445 W/o  $\mathcal{L}^{\text{Time}}$  and  $\mathcal{L}^{\text{Freq}}$  are 0.859 and 0.547 on FD-A dataset. One can conclude that the time- and 446 frequency-based contrastive modules have almost the same contribution to the whole TIC model. 447 The variant W/o  $\mathcal{L}^{\text{Inte}}$  has the lowest ARI on all the datasets. This is because every positive pair 448 consists only of the sample and its augmentation after removing the  $\mathcal{L}^{\text{Inte}}$ , i.e., losing the basic 449 clustering information provided by pseudo-labels. In this case, the class collision issue seriously 450 degrades the clustering performance as discussed in Remark 2. In Fig.4(b), we further show the 451 cosine distances among samples belonging to different clusters. One can see that considering the interactive loss  $\mathcal{L}^{\text{Inte}}$  (with  $\mathcal{L}^{\text{Inte}}$  in Fig.4(b)) can increase the inter-cluster cosine distance with a 452 remarkable gap of 20.4% (i.e., from 1.965 to 2.365 on ECG dataset). It shows that minimizing the 453  $\mathcal{L}^{\text{Inte}}$  indeed increases the distinctiveness of learned embeddings and bring forward better cluster 454 performance. 455



Figure 4: Cosine distance (a) between time- and frequency-based embeddings of the same sample considering the combination (With comb, i.e., full TIC model) and without the combination (W/o Comb) of results from time and frequency domains. The lower distance denotes better encoder learning, because  $\mathbf{g}_i^{\text{time}}$  and  $\mathbf{g}_i^{\text{time}}$  are closer in embedding space. Cosine distance (b) among the concatenated embeddings  $[\mathbf{g}_i^{\text{time}}; \mathbf{g}_i^{\text{time}}]$  of samples belonging to different clusters. The larger intercluster distance represents a better clustering result.

471

456

457

458

459

460

461

462

464

435

**RQ3.** Embedding evaluation. We evaluate the embedding learned by TIC model by performing 472 two other downstream tasks: classification and anomaly detection. The results w.r.t. anomaly detec-473 tion are shown in Appendix A.8. We follow the same protocol as (Franceschi et al., 2019; Tonek-474 aboni et al., 2022), where a multi-class SVM with RBF kernel and a linear classifier are trained on 475 top of the embeddings learned by different models. The 5-fold cross-validation is adopted to train 476 the SVM and the linear classifier. In TIC, the time- and frequency-based embeddings are concate-477 nated  $[\mathbf{g}^{\text{Time}}; \mathbf{g}^{\text{Freq}}]$ . Beyond the aforementioned unsupervised representation learning methods, we 478 consider a K-nearest neighbor classifier (K = 5) equipped with DTW metric (KNN<sub>DTW</sub>) as another 479 baseline. The evaluation results are summarized in Table 4, where the results w.r.t the linear classifier 480 and SVM are shown in the first and second sub-tables. As can be seen, TIC achieves the best perfor-481 mance (colored blue) in 13 out of 18 cases and the second-best performance in another 5 cases. In 482 particular, TIC shows the highest accuracy of 0.9654 when training an SVM classifier, which yields a margin of 4.3% over the best baseline activity2vec (0.9257). It shows that TIC adequately lever-483 ages the information from time and frequency domains to provide more fine-grained embeddings 484 for discriminant. Besides, almost all the unsupervised representation learning baselines have higher 485 accuracy than KNN<sub>DTW</sub>, illustrating the dominance of neural networks in representation learning.

489										
400	Linear	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
490	KNN <sub>DTW</sub>	$.8452 \pm .089$	.6712±.053	.7152±.065	$.5746 \pm .031$	$.6784 \pm .053$	$.3462 \pm .078$	.6745±.043	.8756±.074	$.7127 \pm .054$
491	activity2vec	$.9587 \pm .011$	$.7853 \pm .036$	$.9258 \pm .073$	$.6547 \pm .026$	<u>.<b>8964</b></u> ±.039	$.6578 \pm .038$	$.8941 \pm .085$	$.9245 \pm .004$	<u>.9123</u> ±.030
	BTSF	$.9578 \pm .057$	$.8514 \pm .034$	<u>.9463</u> ±.068	<u>.8637</u> ±.079	$.8875 \pm .068$	$.7123 \pm .029$	$.8745 \pm .075$	$.9521 \pm .056$	$.8654 \pm .028$
492	CLOCS	$.8924 \pm .049$	$.7546 \pm .056$	$.8954 \pm .038$	$.4852 \pm .023$	$.8456 \pm .087$	$.7214 \pm .036$	$.8875 \pm .074$	$.9485 \pm .029$	$.8745 \pm .054$
/03	MHCCL	$.9643 \pm .022$	$.7765 \pm .031$	$.9164 \pm .034$	$.7756 \pm .055$	$.8382 \pm .064$	$.7936 \pm .051$	.9145 <sub>±.029</sub>	$.9654 \pm .005$	$.7363 \pm .098$
455	TFC	<u>.9854</u> ±.007	<b>.9087</b> ±.045	$.9245 \pm .054$	$.7955 \pm .024$	$.8657 \pm .048$	<u>.8834</u> ±.038	$.8921 \pm .051$	$.9478 \pm .025$	$.7982 \pm .069$
494	Triplet	$.9338 \pm .077$	$.8937 \pm .064$	$.9054 \pm .029$	$.6953 \pm .044$	$.8542 \pm .055$	$.8377 \pm .047$	$.8999 \pm .035$	<u>.9785</u> ±.009	$.8032 \pm .013$
405	TS2Vec	$.8687 \pm .035$	$.8546 \pm .054$	$.4887 \pm .067$	$.8365 \pm .081$	$.8922 \pm .039$	$.8643 \pm .059$	$.8456 \pm .062$	$.9087 \pm .029$	$.8266 \pm .002$
495	TSTCC	$.9485 \pm .014$	$.7481 \pm .011$	$.8804 \pm .025$	$.7685 \pm .024$	$.8547 \pm .039$	$.8598 \pm .011$	$.8300 \pm .007$	$.9158 \pm .009$	$.7591 \pm .003$
496	TIC (ours)	<u>.9785</u> ±.057	<u>.9132</u> ±.056	<u>.9571</u> ±.084	<u>.8795</u> ±.067	<u>.9215</u> ±.069	<u>.8965</u> ±.027	<u>.9013</u> ±.039	<u>.9874</u> ±.055	<u>.9251</u> ±.036
-100	SVM	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
497	KNNDTW	$.8452 \pm .089$	$.6712 \pm .053$	$.7152 \pm .065$	$.5746 \pm .031$	$.6784 \pm .053$	$.3462 \pm .078$	$.6745 \pm .043$	$.8756 \pm .074$	$.7127 \pm .054$
400	activity2vec	$.9651 \pm .035$	$.8324 \pm .067$	$.9257 \pm .054$	$.68/1 \pm .036$	$.9031 \pm .089$	$.6982 \pm .056$	$.8874 \pm .074$	$.9483 \pm .036$	$.9245 \pm .057$
498	BISF	$.9687 \pm .045$	$.8789 \pm .037$	$.9128 \pm .061$	$\frac{.8843}{.039} \pm .039$	$.8992 \pm .066$	$.7536 \pm .052$	$.8905 \pm .037$	$.9569 \pm .045$	$.8763 \pm .079$
499	CLOCS	$.9214 \pm .066$	$./895 \pm .074$	$.9214 \pm .056$	$.5123 \pm .081$	$.851/\pm.036$	$./541 \pm .063$	$.890/\pm.046$	$\frac{.9681}{.039}$ ± .039	$.8852 \pm .028$
100	MHCCL	$.9095 \pm .035$	$.7854 \pm .076$	$.9237 \pm .046$	./980±.048	.8/30±.078	$.8005 \pm .031$	.8943±.032	$.9535 \pm .067$	$.7629 \pm .054$
500	TFC	<u>.9743</u> ±.085	<u>.9316</u> ±.037	$.9158 \pm .045$	$.8214 \pm .061$	$.8702 \pm .047$	$\frac{.9317}{\pm .036}$	$.8795 \pm .026$	$.9538 \pm .054$	$.8369 \pm .091$
501	Triplet	$.9502 \pm .045$	$.9125 \pm .067$	$.9056 \pm .057$	$.8506 \pm .075$	$.8722 \pm .049$	$.869/\pm.094$	$.890/\pm.056$	$.9654 \pm .048$	$.8537 \pm .033$
100	TS2Vec	$.8794 \pm .024$	$.8874 \pm .079$	$.5638 \pm .047$	$.8574 \pm .067$	$\frac{.9201}{.043}$	$.8932 \pm .046$	$.8645 \pm .033$	$.9268 \pm .076$	$.8437 \pm .070$
502	TIC (ours)	$.9542 \pm .036$	$.7896 \pm .033$	$\frac{.9214}{9654}$ .096	$.7982 \pm .019$	$.8964 \pm .056$	$.8009 \pm .087$	$.8514 \pm .076$	$9235 \pm .044$	$9391 \pm 021$
503	The (ours)	1000±.005	.049	<u></u>	.081	.004	.074	.026	<u>0740</u> ±.037	<u></u>
				• · ·				DI 141 1144		
504	1	(a) Al	RI of Gesture	dataset		1 -	(b) A	RI with diffe	erent t	
505		anna		Karan A		t=10	) nout update	т		τ
506	0.8		্র	1		0.9 - t=20	2			
507	0.6				~	t=40	, ⊥ T			

486 Table 4: Accuracy (mean±std.deviation) of different methods on benchmark datasets. The results w.r.t 487 the linear classifier and the multi-class SVM are shown in the first and second sub-tables. KNN<sub>DTW</sub> 488 denotes a K-nearest neighbor classifier equipped with DTW metric.



Figure 5: ARI in every epoch with updating pseudo-labels every 40 epochs and without update (a). ARI decreases slightly each time the pseudo-labels are updated by k-means at  $40^{th}$  and  $80^{th}$  epochs (shown by black circles). ARI with different t is shown in (b). The frequent updating of pseudolabels (t = 10) degrades the clustering performance of TIC model.

**RQ4.** Sensitivity of the hyperparameter t. As shown in Fig.1, TIC model updates the pseudo-518 labels generated by k-means every t epochs. In Fig.5(a), we show the ARI of TIC in every epoch 519 with t = 40 and without update for the Gesture dataset. It can be seen that the clustering results 520 generated by k-means keep being improved (ARI=0.347 with epoch  $\in$  [1,39], ARI=0.623 with 521  $epoch \in [40, 79]$ , ARI=0.714 with  $epoch \in [80, 99]$ ) because the embedding learning is constantly 522 optimized. This leads to basic clustering information of higher quality being considered in the 523 interaction loss. Although the ARI of TIC decreases temporarily at each update of the pseudo-label 524 (epoch = 40 and epoch = 80), the final ARI with t = 40 is higher than the one without updates. 525 Fig.5(b) shows the ARI of TIC model with different t. Frequent updating (t = 10) of pseudo-526 labels degrades the performance of TIC, even making it lower than the case of no updating, as the embedding learning is destabilized. we also provide the visualization of the clustering results and 527 the effects of various data augmentation techniques in Appendix A.9 and A.10. 528

529 530

531

508

509

510

511 512

513

514

515

516 517

#### 5 CONCLUSION

532 This paper proposes a trusted and interactive clustering model for time-series data, named TIC, 533 leveraging evidence theory to combine time- and frequency-based information. TIC optimizes the 534 contrastive loss from time and frequency domains, and an interactive loss calculated based on the 535 pseudo-labels. Uncertainty in time- and frequency-based clustering results are quantified by mass 536 functions that are combined by Dempster's rule to produce the trusted clustering results. Experimen-537 tal results show the superior clustering performance of TIC brought by the combination of clustering results from time and frequency domains, as well as the consideration of interactive loss. The em-538 bedding learned by TIC is also shown to perform well on the classification and anomaly detection tasks.

# 540 REFERENCES

549

550

551

552

565

566

567

- Karan Aggarwal, Shafiq Joty, Luis Fernandez-Luque, and Jaideep Srivastava. Adversarial unsupervised representation learning for activity time-series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 834–841, 2019.
- 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.
  - Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, Jorge Luis Reyes-Ortiz, et al. A public domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, pp. 3, 2013.
- Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi.
   A theoretical analysis of contrastive unsupervised representation learning. In *36th International Conference on Machine Learning, ICML 2019*, pp. 9904–9923. International Machine Learning Society (IMLS), 2019.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
  - <sup>0</sup> E Oran Brigham. *The fast Fourier transform and its applications*. Prentice-Hall, Inc., 1988.
- Ricardo JGB Campello, Davoud Moulavi, and Jörg Sander. Density-based clustering based on hierarchical density estimates. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 160–172. Springer, 2013.
  - Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3558–3568, 2021.
- Bertrand Charpentier, Daniel Zügner, and Stephan Günnemann. Posterior network: Uncertainty estimation without ood samples via density-based pseudo-counts. Advances in Neural Information Processing Systems, 33:1356–1367, 2020.
- 572 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
  573 contrastive learning of visual representations. In *International Conference on Machine Learning*,
  574 pp. 1597–1607. PMLR, 2020.
- 575
   576
   576
   576
   576
   577
   578
   Yanping Chen, Eamonn Keogh, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, and Gustavo Batista. The ucr time series classification archive, July 2015. www.cs.ucr.edu/ ~eamonn/time\_series\_data/.
- Gari D Clifford, Chengyu Liu, Benjamin Moody, H Lehman Li-wei, Ikaro Silva, Qiao Li, AE Johnson, and Roger G Mark. Af classification from a short single lead ecg recording: The physionet/computing in cardiology challenge 2017. In 2017 Computing in Cardiology (CinC), pp. 1–4. IEEE, 2017.
- Leon Cohen. *Time-frequency analysis*, volume 778. Prentice hall New Jersey, 1995.
- AP DEMPSTER. Upper and lower probabilities induced by a multivalued mapping. Annals of Mathematical Statistics, 38:325–339, 1967.
- Thierry Denœux. Conjunctive and disjunctive combination of belief functions induced by nondistinct bodies of evidence. *Artificial Intelligence*, 172(2-3):234–264, 2008.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
   bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 593 Yuan-Wei Du and Jiao-Jiao Zhong. Generalized combination rule for evidential reasoning approach and dempster–shafer theory of evidence. *Information Sciences*, 547:1201–1232, 2021.

- 594 Jiangyong Duan and Lili Guo. Variable-length subsequence clustering in time series. IEEE Trans-595 actions on Knowledge and Data Engineering, 34(2):983–995, 2023. 596 Didler Dubois and Henri Prade. Representation and combination of uncertainty with belief functions 597 and possibility measures. *Computational intelligence*, 4(3):244–264, 1988. 598 Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee Keong Kwoh, Xiaoli Li, and 600 Cuntai Guan. Time-series representation learning via temporal and contextual contrasting. arXiv 601 preprint arXiv:2106.14112, 2021. 602 603 Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee-Keong Kwoh, Xiaoli Li, and Cuntai Guan. Self-supervised contrastive representation learning for semi-supervised time-series 604 classification. arXiv preprint arXiv:2208.06616, 2022. 605 606 Leonardo N Ferreira and Liang Zhao. Time series clustering via community detection in networks. 607 Information Sciences, 326:227–242, 2016. 608 609 Patrick Flandrin. *Time-frequency/time-scale analysis*. Academic press, 1998. 610 Mihai Cristian Florea, Anne-Laure Jousselme, Éloi Bossé, and Dominic Grenier. Robust combina-611 tion rules for evidence theory. Information Fusion, 10(2):183–197, 2009. 612 613 V Fortuin, M Hüser, F Locatello, H Strathmann, and G Rätsch. Som-vae: Interpretable discrete 614 representation learning on time series. In International Conference on Learning Representations 615 (ICLR 2020), 2020. 616 617 Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised scalable representation learning for multivariate time series. Advances in neural information processing systems, 32, 618 2019. 619 620 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model 621 uncertainty in deep learning. In International Conference on Machine Learning, pp. 1050–1059. 622 PMLR, 2016. 623 624 Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G 625 Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic 626 signals. circulation, 101(23):e215-e220, 2000. 627 628 Chaoyu Gong, Yongbin Li, Di Fu, Yong Liu, Pei-hong Wang, and Yang You. Self-reconstructive 629 evidential clustering for high-dimensional data. In 2022 IEEE 38th International Conference on 630 *Data Engineering (ICDE)*, pp. 2099–2112. IEEE, 2022. 631 632 D Guijo-Rubio, AM Duran-Rosal, PA Gutierrez, A Troncoso, and C Hervas-Martinez. Time-series clustering based on the characterization of segment typologies. IEEE Transactions on Cybernet-633 ics, 51(11):5409-5422, 2021. 634 635 Zongbo Han, Changqing Zhang, Huazhu Fu, and Joey Tianyi Zhou. Trusted multi-view classifica-636 tion. arXiv preprint arXiv:2102.02051, 2021. 637 638 Zongbo Han, Changqing Zhang, Huazhu Fu, and Joey Tianyi Zhou. Trusted multi-view classi-639 fication with dynamic evidential fusion. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(2):2551-2566, 2022. 640 641 Toshitaka Hayashi, Dalibor Cimr, Filip Studnička, Hamido Fujita, Damián Bušovský, Richard Cim-642 ler, and Ali Selamat. Distance-based one-class time-series classification approach using local 643 cluster balance. Expert Systems with Applications, 235:121201, 2024. 644 645 Rob J Hyndman and George Athanasopoulos. Forecasting: principles and practice. OTexts, 2018. 646
- 647 Fumitada Itakura. Minimum prediction residual principle applied to speech recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 23(1):67–72, 1975.

659

660

661

685

688

689

- 648 Audun Jsang. Subjective Logic: A formalism for reasoning under uncertainty. Springer Publishing 649 Company, Incorporated, 2018. 650
- Bob Kemp, Aeilko H Zwinderman, Bert Tuk, Hilbert AC Kamphuisen, and Josefien JL Oberye. 651 Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the eeg. 652 *IEEE Transactions on Biomedical Engineering*, 47(9):1185–1194, 2000. 653
- 654 Dani Kiyasseh, Tingting Zhu, and David A Clifton. Clocs: Contrastive learning of cardiac signals 655 across space, time, and patients. In International Conference on Machine Learning, pp. 5606-5615. PMLR, 2021. 656
- Anna-Kathrin Kopetzki, Bertrand Charpentier, Daniel Zügner, Sandhya Giri, and Stephan 658 Günnemann. Evaluating robustness of predictive uncertainty estimation: Are dirichlet-based models reliable? In International Conference on Machine Learning, pp. 5707–5718. PMLR, 2021.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive 662 uncertainty estimation using deep ensembles. Advances in neural information processing systems, 663 30, 2017. 664
- 665 Christian Lessmeier, James Kuria Kimotho, Detmar Zimmer, and Walter Sextro. Condition mon-666 itoring of bearing damage in electromechanical drive systems by using motor current signals of 667 electric motors: A benchmark data set for data-driven classification. In PHM Society European 668 Conference, volume 3, 2016.
- 669 Jiayang Liu, Lin Zhong, Jehan Wickramasuriya, and Venu Vasudevan. uwave: Accelerometer-670 based personalized gesture recognition and its applications. Pervasive and Mobile Computing, 5 671 (6):657-675, 2009. 672
- 673 Xinwang Liu, Xinzhong Zhu, Miaomiao Li, Lei Wang, Chang Tang, Jianping Yin, Dinggang Shen, Huaimin Wang, and Wen Gao. Late fusion incomplete multi-view clustering. IEEE Transactions 674 on Pattern Analysis and Machine Intelligence, 41(10):2410-2423, 2018. 675
- 676 Q Ma, S Li, W Zhuang, J Wang, and D Zeng. Self-supervised time series clustering with model-677 based dynamics. IEEE Transactions on Neural Networks and Learning Systems, 32(9):3942-678 3955, 2022. 679
- Qianli Ma, Sen Li, Lifeng Shen, Jiabing Wang, Jia Wei, Zhiwen Yu, and Garrison W Cottrell. End-680 to-end incomplete time-series modeling from linear memory of latent variables. IEEE Transac-681 tions on Cybernetics, 50(12):4908-4920, 2019a. 682
- 683 Qianli Ma, Jiawei Zheng, Sen Li, and Gary W Cottrell. Learning representations for time series 684 clustering. Advances in neural information processing systems, 32, 2019b.
- Wenjun Ma, Yuncheng Jiang, and Xudong Luo. A flexible rule for evidential combination in 686 dempster-shafer theory of evidence. Applied Soft Computing, 85:105512, 2019c. 687
  - Andrey Malinin, Bruno Mlodozeniec, and Mark Gales. Ensemble distribution distillation. In International Conference on Learning Representations, 2019.
- Qianwen Meng, Hangwei Qian, Yong Liu, Lizhen Cui, Yonghui Xu, and Zhiqi Shen. Mhccl: 691 Masked hierarchical cluster-wise contrastive learning for multivariate time series. In Proceed-692 ings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 9153-9161, 2023. 693
- 694 Matthew Middlehurst, Patrick Schäfer, and Anthony Bagnall. Bake off redux: a review and experimental evaluation of recent time series classification algorithms. Data Mining and Knowledge 695 Discovery, pp. 1–74, 2024. 696
- 697 John Paparrizos and Luis Gravano. k-shape: Efficient and accurate clustering of time series. In Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, pp. 699 1855-1870, 2015. 700
- Clément Péalat, Guillaume Bouleux, and Vincent Cheutet. Improved time series clustering based 701 on new geometric frameworks. Pattern Recognition, 124:108423, 2023.

702 François Petitjean, Alain Ketterlin, and Pierre Gançarski. A global averaging method for dynamic 703 time warping, with applications to clustering. Pattern Recognition, 44(3):678–693, 2011. 704 William M Rand. Objective criteria for the evaluation of clustering methods. Journal of the Ameri-705 can Statistical association, 66(336):846-850, 1971. 706 707 Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classifica-708 tion uncertainty. Advances in neural information processing systems, 31, 2018. 709 G. Shafer. A mathematical theory of evidence. Princeton university press, 1976. 710 711 Ryan Soklaski, Michael Yee, and Theodoros Tsiligkaridis. Fourier-based augmentations for im-712 proved robustness and uncertainty calibration. arXiv preprint arXiv:2202.12412, 2022. 713 Zhi-gang Su, Thierry Denoeux, Yong-sheng Hao, and Ming Zhao. Evidential k-nn classification 714 with enhanced performance via optimizing a class of parametric conjunctive t-rules. Knowledge-715 Based Systems, 142:7-16, 2018. 716 717 Chi Ian Tang, Ignacio Perez-Pozuelo, Dimitris Spathis, and Cecilia Mascolo. Exploring contrastive learning in human activity recognition for healthcare. arXiv preprint arXiv:2011.11542, 2020. 718 719 Yongqiang Tang, Yuan Xie, Xuebing Yang, Jinghao Niu, and Wensheng Zhang. Tensor multi-elastic 720 kernel self-paced learning for time series clustering. IEEE Transactions on Knowledge & Data 721 *Engineering*, 33(03):1223–1237, 2021. 722 Sana Tonekaboni, Danny Eytan, and Anna Goldenberg. Unsupervised representation learning for 723 time series with temporal neighborhood coding. arXiv preprint arXiv:2106.00750, 2022. 724 725 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of Machine 726 Learning Research, 9(11), 2008. 727 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 728 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa-729 tion processing systems, 30, 2017. 730 731 Siwei Wang, Xinwang Liu, En Zhu, Chang Tang, Jiyuan Liu, Jingtao Hu, Jingyuan Xia, and Jian-732 ping Yin. Multi-view clustering via late fusion alignment maximization. In International Joint 733 Conference on Artificial Intelligence, pp. 3778–3784, 2019. 734 Siwei Wang, Xinwang Liu, Li Liu, Sihang Zhou, and En Zhu. Late fusion multiple kernel clustering 735 with proxy graph refinement. IEEE Transactions on Neural Networks and Learning Systems, 736 2021. 737 Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Cost: Contrastive learn-738 ing of disentangled seasonal-trend representations for time series forecasting. In International 739 Conference on Learning Representations, 2021. 740 741 Ronald R Yager. On the dempster-shafer framework and new combination rules. Information sci-742 ences, 41(2):93-137, 1987. 743 Jaewon Yang and Jure Leskovec. Patterns of temporal variation in online media. In Proceedings of 744 the fourth ACM International Conference on Web Search and Data Mining, pp. 177–186, 2011. 745 746 Ling Yang and Shenda Hong. Unsupervised time-series representation learning with iterative bi-747 linear temporal-spectral fusion. In International Conference on Machine Learning, pp. 25038-748 25054. PMLR, 2022. 749 Zhihan Yue, Yujing Wang, Juanyong Duan, Tianmeng Yang, Congrui Huang, Yunhai Tong, and 750 Bixiong Xu. Ts2vec: Towards universal representation of time series. In Proceedings of the AAAI 751 Conference on Artificial Intelligence, volume 36, pp. 8980–8987, 2022. 752 753 Kexin Zhang, Qingsong Wen, Chaoli Zhang, Rongyao Cai, Ming Jin, Yong Liu, James Y Zhang, Yuxuan Liang, Guansong Pang, Dongjin Song, et al. Self-supervised learning for time series anal-754 ysis: Taxonomy, progress, and prospects. IEEE Transactions on Pattern Analysis and Machine 755

Intelligence, 2024.

Qin Zhang, Jia Wu, Peng Zhang, Guodong Long, and Chengqi Zhang. Salient subsequence learning for time series clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41 (09):2193–2207, 2019.

 Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. Self-supervised contrastive pre-training for time series via time-frequency consistency. *Advances in Neural Information Processing Systems*, 35:3988–4003, 2022.

763 764

765

- A APPENDIX
- 766 767 A.1 Related work

We review time-series clustering methods, contrastive learning methods toward time-series data and uncertain learning methods, respectively.

771 A.1.1 TIME-SERIES CLUSTERING

Time-series clustering methods can be roughly classified into two families: raw-data-based methodsand feature-based methods.

Raw-data-based methods. These methods primarily modify the distance metric to accommodate the specific characteristics of time-series data. K-DBA (Petitjean et al., 2011) combines *k*-means and dynamic time warping (DTW) (Itakura, 1975) to achieve improved alignment. To reveal the temporal dynamics, K-SC method (Yang & Leskovec, 2011) utilizes a similarity metric that is invariant to scaling and shifting. K-Shape (Paparrizos & Gravano, 2015) considers the shapes of time-series by employing a normalized cross-correlation metric. The mentioned methods often exhibit sensitivity to outliers and noise because they consider all time-step points (Ma et al., 2019b).

782 **Feature-based methods.** Feature-based methods typically involve two stages, where the input time-783 series samples are transformed into informative features first and clustering algorithms are then con-784 ducted on these features (Tang et al., 2021; Fortuin et al., 2020). In conjunction with k-means, TNC 785 (Tonekaboni et al., 2022) ensures the distribution of signals from the neighborhood is distinguishable from the distribution of non-neighboring signals. Authors in (Tang et al., 2021) map the raw 786 time-series space into multiple kernel spaces via elastic distance measure functions and resort to 787 a self-paced learning paradigm to group time-series samples. Authors in (Péalat et al., 2023) em-788 bed the time-series onto the Stiefel manifold to obtain the geometric representations of time-series 789 samples. STCN (Ma et al., 2022) optimizes the feature extraction and clustering simultaneously, 790 through a recurrent neural network and a self-supervised clustering module. However, almost all the 791 feature-based methods do not adopt an information fusion perspective to incorporate both time and 792 frequency domain information for the purpose of clustering.

793 794 795

#### A.1.2 CONTRASTIVE LEARNING TOWARD TIME-SERIES DATA

Contrastive learning is a well-known form of self-supervised learning and aims to train an encoder 796 that maps original inputs into an embedding space. The objective is to bring positive sample pairs 797 closer together, while pushing negative sample pairs (comprising the original augmentation and 798 an alternative augmentation of a different input sample) apart (Chen et al., 2020). Compared to 799 CV (Changpinyo et al., 2021) and NLP (Devlin et al., 2018), contrastive learning in the context of 800 time-series data has been less explored, mainly due to the difficulty of capturing crucial invariance 801 properties specific to time-series. TF-C (Zhang et al., 2022) expects that time-based and frequency-802 based representations of the same sample are located close together in the time-frequency space, and 803 embeds the time-based neighborhood of a sample close to its frequency-based neighborhood. BTSF 804 (Yang & Hong, 2022) utilizes sample-level augmentation with a dropout on a time-series sample, 805 and devises the iterative bilinear temporal-spectral fusion to generate discriminative embeddings. 806 CoST (Woo et al., 2021) comprises both time domain and frequency domain contrastive losses to 807 learn seasonal representations for long sequence time-series forecasting. In addition to these three methods involving both time and frequency domains, other methods mainly focus on the augmen-808 tations implemented in time domain, such as transformation invariance (e.g., SimCLR (Tang et al., 2020; Chen et al., 2020)) and contextual invariance (e.g., TS2vec (Yue et al., 2022) and TS-TCC

810 (Eldele et al., 2021)). In previous works, the loss information from the time- and frequency-domain 811 is captured in a compositional way, e.g., TF-C simply sums the loss functions from the two domains 812 and the consistency loss function, and BTSF solely implements the data augmentation in the time do-813 main. To the best of our knowledge, this work is the first one that directly combines time-frequency 814 domain information to leverage information fusion for time-series clustering.

#### 816 A.1.3 UNCERTAINTY-BASED LEARNING

817 Some efforts (Gal & Ghahramani, 2016; Lakshminarayanan et al., 2017; Charpentier et al., 2020) 818 have been made to enable the neural network to estimate the uncertainty of the output. Evidential 819 network (Sensoy et al., 2018) incorporates subjective logic to model the Dirichlet distribution. Post-820 Net (Malinin et al., 2019) utilizes normalizing flow and Bayesian loss during training to estimate 821 uncertainty. TMC (Han et al., 2021) and ETMC (Han et al., 2022) introduce a variational Dirich-822 let distribution to characterize the distribution of the class probabilities in multi-view classification. 823 Ensemble distribution distillation (Malinin et al., 2019) leverages the predictions of multiple mod-824 els to estimate the uncertainty. Authors in (Kopetzki et al., 2021) apply median smoothing to the 825 Dirichlet model and enhance the capability of the model to handle adversarial examples. Unlike the above methods, our method is perhaps the first attempt to estimate the uncertainty of the outputs 826 from a neural network oriented to the clustering problem. Further, we fuse the output with uncer-827 tainty estimation from the time and frequency domains to obtain trusted clustering results under the 828 framework of evidence theory. 829

#### A.2 PROOFS OF PROPOSITIONS 1 AND 2

**Proposition 1:** A large  $m_{\Omega}^{T}$  does not lead to a large  $m_{\Omega}^{Co}$ , when one of  $m_{c}^{F}$  is large and  $m_{\Omega}^{F}$  is small. In particular,  $\mathbf{m}^{Co}$  is identical to  $\mathbf{m}^{T}$ , if  $\mathbf{m}^{F}$  is totally uncertain (i.e.,  $m_{\Omega}^{F} = 1$ ).

Proof:

815

830

831 832

833 834 835

842

845 846

854 855

856

$$\begin{split} m_c^{\mathrm{Co}} &= \frac{m_c^{\mathrm{T}} m_{\Omega}^{\mathrm{F}} + m_c^{\mathrm{F}} m_{\Omega}^{\mathrm{T}} + m_c^{\mathrm{T}} m_c^{\mathrm{F}}}{m_{\Omega}^{\mathrm{F}} m_{\Omega}^{\mathrm{T}} + \sum_{v=1}^{C} m_v^{\mathrm{T}} m_v^{\mathrm{F}} + (1 - m_{\Omega}^{\mathrm{T}}) m_{\Omega}^{\mathrm{F}} + (1 - m_{\Omega}^{\mathrm{F}}) m_{\Omega}^{\mathrm{T}}} \\ &= \frac{m_c^{\mathrm{T}} m_{\Omega}^{\mathrm{F}} + m_c^{\mathrm{F}} m_{\Omega}^{\mathrm{T}} + m_c^{\mathrm{T}} m_c^{\mathrm{F}}}{\sum_{v=1}^{C} m_v^{\mathrm{T}} m_v^{\mathrm{F}} + m_{\Omega}^{\mathrm{T}} + m_{\Omega}^{\mathrm{F}} - m_{\Omega}^{\mathrm{F}} m_{\Omega}^{\mathrm{T}}} \end{split}$$

Considering the worst case with  $m_{\Omega}^{\rm F} = 1, m_v^{\rm F} = 0, v = 1, 2, \cdots, n$ , then we can get  $m_c^{\rm Co} = m_c^{\rm T}$ . 843 844

**Proposition 2:** The  $m_{\Omega}^{\text{Co}}$  is monotonically increasing with  $m_{\Omega}^{\text{T}}$  and  $m_{\Omega}^{\text{F}}$ .

Proof:

$$m_{\Omega}^{\text{Co}} = \frac{m_{\Omega}^{\text{T}} m_{\Omega}^{\text{F}}}{\sum_{v=1}^{C} (m_{v}^{\text{T}} m_{v}^{\text{F}} + m_{v}^{\text{T}} m_{\Omega}^{\text{F}} + m_{v}^{\text{F}} m_{\Omega}^{\text{T}}) + m_{\Omega}^{\text{T}} m_{\Omega}^{\text{F}}}$$
$$= \frac{1}{\sum_{v=1}^{C} (\frac{m_{v}^{\text{T}} m_{v}^{\text{F}}}{m_{\Omega}^{\text{T}} m_{\Omega}^{\text{F}}} + \frac{m_{v}^{\text{T}}}{m_{\Omega}^{\text{T}}} + \frac{m_{v}^{\text{F}}}{m_{\Omega}^{\text{F}}}) + 1}$$

As can be seen,  $m_{\Omega}^{\text{Co}}$  increases as  $m_{\Omega}^{\text{T}}$  and  $m_{\Omega}^{\text{F}}$  increase.

#### A.3 DESCRIPTION OF DATASETS

ECG (Clifford et al., 2017) is from the 2017 PhysioNet Challenge focusing on classifying ECG 858 recordings, where single-lead ECG measures four different underlying conditions of cardiac ar-859 rhythmias. EMG (Goldberger et al., 2000) consists of single-channel Electromyograms (EMGs) 860 recorded from the tibialis anterior muscle of three healthy volunteers suffering from myopathy and 861 neuropathy. HAR (Anguita et al., 2013) contains recordings of 30 health volunteers performing daily activities, including walking, walking upstairs, walking downstairs, sitting, standing, and ly-862 ing. Gesture (Liu et al., 2009) contains accelerometer measurements of eight simple gestures that 863 differ based on the paths of hand movement. FD-A/FD-B (Lessmeier et al., 2016) corresponds to

Faulty Detection Condition A (FD-A) and Faulty Detection Condition B (FD-B), which are generated by an electromechanical drive system that monitors the condition of rolling bearings and detects their failures. SleepEEG (Kemp et al., 2000) are collected from 82 healthy subjects and contains 153 whole-night sleeping electroencephalography (EEG) recordings. Epilepsy (Andrzejak et al., 2001) contains single-channel EEG measurements from 500 subjects and the corresponding brain activity is recorded for 23.6 seconds. The details about the UCR datasets can be found in (Chen et al., 2015).

871 872

873

874

875 876

877

878

879

880

882

883

884

885

887

889

890 891

894

895

896

897

899

900

901

902

903 904

905

906

907 908

909

910

A.4 DESCRIPTION OF BASELINES

**Baselines.** To answer Q1, we consider the following 10 time-series clustering methods in the comparison experiment.

- **DTCR** (Ma et al., 2019b): it integrates the temporal reconstruction and *k*-means objective into the seq2seq model. By proposing a fake-sample generation strategy and auxiliary classification task, it can learn cluster-specific temporal representations.
- **k-shape** (Paparrizos & Gravano, 2015): it relies on a scalable iterative refinement procedure and uses a normalized cross-correlation measure to consider the shapes of time-series.
- **SOM-VAE** (Fortuin et al., 2020): it overcomes the non-differentiability in discrete representation learning and presents a gradient-based version of the traditional self-organizing map (SOM) algorithm.
- **STCN** (Ma et al., 2022): it optimizes the feature extraction and clustering simultaneously, where an RNN conducts the reconstruction of time-series and a self-supervised module to obtain the clustering result.
- **TMEK** (Tang et al., 2021): it maps the raw time-series space into multiple kernel spaces via elastic distance measure functions, and resorts to the tensor-constraint-based self-representation subspace clustering approach.
- TNC (Tonekaboni et al., 2022): it uses the local smoothness of a signal's generative process to define neighborhoods in time-series. By using a de-biased contrastive objective, it learns time-series representations that are input to *k*-means to produce the clusteri.ng results.
  - $\mathbf{TS3C}_{ch}$  (Guijo-Rubio et al., 2021): it consists of two stages, where a least squares polynomial technique is first used to segment the time-series and the hierarchical clustering is applied to all the mapped segmentations.
  - UMAP (Péalat et al., 2023): it embeds the time-series onto higher-dimensional spaces and is in conjunction with HDBSCAN algorithm (Campello et al., 2013) to obtain the results.
  - USSL (Zhang et al., 2019): it integrates the shapelet learning, shapelet regularization, spectral analysis and pseudo-label to automatically learn shapelets to group time-series samples.
  - VLSC (Duan & Guo, 2023): it minimizes the inner time-series clustering error under time-series cover constraints, where the time-series lengths can be variable.

To answer Q3, we evaluate the embeddings learned by TIC on the other downstream tasks (i.e., classification and anomaly detection). The included baselines are 8 unsupervised representation learning methods for time-series.

- **activity2vec** (Aggarwal et al., 2019): it learns the representations with three components, the cooccurrence and magnitude of the activity levels in a time-series sample, neighboring context of the time-series, and promoting subject-invariance with adversarial training.
- BTSF (Yang & Hong, 2022): it utilizes the sample-level augmentation with a simple dropout on the time-series dataset and devise the iterative bilinear temporal-spectral fusion to encode the affinities of abundant time-frequency pairs.
- CLOCS (Kiyasseh et al., 2021): it encourages the time-series representations across space, time, and patients to be similar to one another for the physiological data.
- MHCCL (Meng et al., 2023): it exploits semantic information obtained from the hierarchical structure consisting of multiple latent partitions for multivariate time-series.

- **TFC** (Zhang et al., 2022): it learns the representations of time-series by positing that embedding a time-based neighborhood of a sample should be close to its frequency-based neighborhood.
- **Triplet** (Franceschi et al., 2019): it learns general-purpose representations by combining an encoder based on causal dilated convolutions with a novel triplet loss employing time-based negative sampling.
- **TS2Vec** (Yue et al., 2022): it performs contrastive learning in a hierarchical way over augmented context views and obtains a contextual representation for each timestamp.
- **TSTCC** (Eldele et al., 2022): it proposes time-series-specific weak and strong augmentations and learns discriminative representations in a contextual contrasting module.

#### A.5 MORE STATISTICAL TEST

We also compare the performance of methods using the Friedman test and Nemenyi test by setting the significant level to 0.05, which is shown in the Fig.6. TIC significantly outperforms the baselines except for SOM-VAE and TS3C<sub>ch</sub> in all cases.



Figure 6: Comparison between TIC and other baselines with the Nemenyi test. The lowest (best) ranks are to the right, and thus methods on the right sides are considered to be better. The groups of baselines that are not significantly different from TIC and are connected by red.

A.6 NMI, ACC AND RUNNING TIME RESULTS

Overall, our TIC model wins 235 and has tied performance on 34 out of 270 trials, when it is statistically compared with 10 baselines based on all the three metrics. Compared to the baselines, TIC consumes comparative runtime but achieves the best clustering accuracy.

#### A.7 MORE ABLATION STUDY

We conduct other ablation studies to evaluate the importance of using the Dempster's rule in the developed TIC model. Concretely, we compare TIC model with the following variants:

- TIC-cau/TIC-bol/TIC-tri: Dempster's rule is replaced by other combination rule, i.e., the cautious conjunctive rule (Denœux, 2008), the bold conjunctive rule (Denœux, 2008), the parametric triangular-norm-based rule (Su et al., 2018), which do not rely on the assumption that the items of evidence are independent of each other;

Table 5: NMI, ACC (mean $\pm$ std.deviation) and running time of different algorithms on benchmark datasets. The •/• indicates whether TIC is statistically superior/inferior to a certain comparing baseline based on the paired *t*-test at a 0.05 significance level. The statistics of win/tie/loss are shown in the last row of each sub-table. The best and the second-best results, between which the performance gaps are shown in the row named "gap" in each sub-table, are colored blue and red.

NMI	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
DTCR	.802±.02•	.745±.01•	.653±.02•	.766±.02●	.862±.01•	<u>.822</u> ±.02°	.492±.02•	.802±.02•	.663±.02•
k-shape	.752±.02●	.591±.01•	.745±.01•	$.734 \pm .02 \bullet$	.836±.02•	$.756 \pm .03$	$.692 \pm .02 \bullet$	.834±.03●	.843±.03●
SOM-VAE	<u>.892</u> ±.03●	.632±.02•	<u>.822</u> ±.02●	.888±.01	.789±.03●	.682±.02•	.722±.02•	.929±.01•	.825±.03●
STCN	$1_{\pm 0}$	.745±.03●	.669±.02●	.793±.01•	<u>.894</u> ±.02	.726±.03●	.833±.02•	.802±.01•	.639±.02●
TMEK	.696±.02●	.653±.02●	.738±.01•	.859±.01●	.721±.02•	.730±.01•	.793±.03●	$.952 \pm .02$	.703±.01•
TNC	.743±.01•	<u>.746</u> ±.02●	.654±.01●	.833±.02•	.751±.02•	.652±.01•	.799±.02●	.854±.01•	.738±.01•
TS3C <sub>ch</sub>	$.752 \pm .01 \bullet$	$.726 \pm .01 \bullet$	.753±.02•	$.871 \pm .02$	.754 <sub>±.04</sub> ●	.752 <sub>±.03</sub> ●	<u>.921</u> ±.01	$.949 \pm .01$	$.874 \pm .01$
UMAP	.863±.01•	$.720 \pm .01 \bullet$	.610±.02•	.827±.02•	.714±.04•	.746±.03●	.840±.01•	<u>.954</u> ±.01	<u>.882</u> ±.01
USSL	.842±.01•	$.702 \pm .01 \bullet$	.496±.02●	$.568 \pm .02 \bullet$	.652±.02•	.705±.03●	.900±.02●	.895±.01•	.736±.02●
VLSC	$\underline{1}_{\pm 0}$	.734±.01•	.578±.02●	.724±.02•	.630±.03●	.339±.02•	.754±.01●	.831±.02•	.731±.01•
TIC (ours)	$\underline{1}_{\pm 0}$	.786±.01	<u>.887</u> ±.01	<u>.898</u> ±.02	<u>.900</u> ±.01	<u>.784</u> ±.02	<u>.942</u> ±.01	<u>.965</u> ±.02	<u>.891</u> ±.01
gap	.108	.040	.065	.010	.006	.038	.021	.011	.009
win/tie/loss	8/2/0	10/0/0	10/0/0	8/2/0	9/1/0	8/1/1	9/1/0	7/3/0	8/2/0
ACC	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
DTCR	.792±.02•	.841±.01•	.621±.02•	.892±.03●	.906±.02•	<u>.938</u> ±.01	.620±.01•	.921±.02•	.657±.02●
k-shape	.715±.02●	.669±.01●	.802±.03●	.856±.02●	<u>.925</u> ±.01•	.837±.02•	.840±.01•	.931±.02•	.831±.03●
SOM-VAE	.885±.03●	.879±.02●	<u>.839</u> ±.01●	<b>.940</b> ±.01	$.862 \pm .02 \bullet$	.781±.02•	.874±.02●	.952 <sub>±.02</sub> ●	$.841 \pm .01 \bullet$
STCN	$\underline{1}_{\pm 0}$	.785±.02●	.805±.02●	.839±.01•	.919±.03●	.846±.02●	.826±.02●	.810±.01•	.605±.02●
TMEK	.742±.02•	.795±.02●	.800±.01•	.863±.01•	.752±.02•	.822±.01•	.793±.03●	.954±.02●	.678±.01●
TNC	.838±.01•	<u>.921</u> ±.02	$.726 \pm .01 \bullet$	.896±.02●	.821±.03•	$.675 \pm .01 \bullet$	.820±.03●	.829±.01•	.771±.01•
$TS3C_{ch}$	.798±.02●	.842±.02●	$.732 \pm .01 \bullet$	.889±.01●	.796±.03●	$.726 \pm .01 \bullet$	$.933 \pm .01$	.950±.02●	<u>.895</u> ±.02
UMAP	$.876 \pm .01 \bullet$	.846±.01●	$.705 \pm .02 \bullet$	.909±.02●	.731±.02•	.863±.01•	$.921 \pm .01 \bullet$	<u>.988</u> ±.03	$.867 \pm .01$
USSL	<u>.932</u> ±.01●	.822±.01●	.734±.02●	.619±.02●	.658±.04●	.639±.03●	.885±.01●	$.978 \pm .01$	.751±.01●
VLSC	$\underline{1}_{\pm 0}$	.846±.01●	.687±.02●	.879±.02●	.786±.03●	.602±.02●	<u>.941</u> ±.01	$.977 \pm .02$	.742±.01•
TIC (ours)	$1_{\pm 0}$	<u>.935</u> ±.01	<u>.869</u> ±.01	<u>.959</u> ±.01	<u>.982</u> ±.01	<u>.950</u> ±.02	<u>.957</u> ±.01	<u>.997</u> ±.01	<u>.914</u> ±.02
gap	.068	.014	.030	.019	.057	.012	.016	.009	.019
win/tie/loss	8/2/0	9/1/0	10/0/0	9/1/0	10/0/0	9/1/0	8/2/0	7/3/0	8/2/0
Time	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR
DTCR	1.65	412.2	113.2	8.5	213.2	324.3	3486.3	143.7	154.3
k-shape	0.32	541.1	132.7	14.3	300.9	452.9	3688.2	168.3	217.8
SOM-VAE	2.51	365.8	101.7	7.2	196.2	331.1	3348.2	121.9	169.8
STCN	0.85	5421.1	2117.3	3.4	1837.6	2913.2	27986.3	2684.9	186.9
TMEK	1.89	432.4	99.4	6.3	145.7	217.3	2987.6	119.3	197.9
TNC	2.07	500.7	145.3	5.2	164.8	213.7	3114.7	192.3	208.9
$TS3C_{ch}$	1.54	472.1	123.6	10.9	151.8	207.9	3864.3	143.5	178.3
UMAP	1.25	625.1	200.6	7.6	275.9	423.9	3004.7	287.9	190.6
USSL	2.09	587.3	132.4	4.7	167.3	246.4	3654.8	169.3	183.7
VLSC	1.74	584.6	241.1	5.6	272.2	384.3	3845.1	300.8	223.7
TIC (ours)	2.09	576.3	135.3	6.8	206.2	312.9	3884.9	249.9	199.8

1000 1001 1002

1003

1004

> • TIC-comc/TIC-GC/TIC-yage/TIC-dubo/TIC-RCR: Dempster's rule is replaced by other combination rule, i.e., the COMC rule (Ma et al., 2019c), the GC rule (Du & Zhong, 2021), Yager's rule (Yager, 1987), Dubois-Prade's rule (Dubois & Prade, 1988), the RCR rule (Florea et al., 2009) which have their own ways to deal with the high-conflicted mass functions.

> > Table 6: ARI values (mean<sub>±std.deviation</sub>) of different variants of TIC.

ARI	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR	Average
TIC×	.845	.729	.711	.921	.895	.803	.832	.893	.873	
TIC+	.839	.698	.687	.909	.910	.821	.850	.869	.867	
TIC-cau	.965	.873	.783	.924	.924	.839	.892	.931	.832	
TIC-bol	.981	.869	.789	.920	.867	.841	.882	.965	.809	
TIC-tri	1	<u>.879</u>	.772	.915	.926	.839	.889	.978	.867	
TIC-comc	.987	.872	.754	.928	.927	.828	<u>.902</u>	.963	.856	
TIC-GC	.979	.871	.768	.916	.913	.819	.856	.952	.873	
TIC-yage	1	.869	.781	.924	.918	.834	.883	<u>.980</u>	.830	
TIC-dubo	.981	.853	.789	.916	.904	.842	.879	.980	.815	
TIC-RCR	1	.874	.792	.904	.919	.836	.882	.958	.871	
TIC (Full model)	1	.879	.796	.931	.939	.853	.896	.980	.883	

1017 1018 Compared to the full model, the average ARI of TIC<sub>×</sub> and TIC<sub>+</sub> decrease by 0.912 - 0.834 = 0.0781019 and 0.912 - 0.833 = 0.079. This suggests that it is more reasonable to use Dempster's rule to fuse 1020 the clustering results from time and frequency domains. The potential reasons for this result are 1021 twofold: (1) as shown in Eq.(9), Dempster's rule includes the  $m_{ic}^{\text{Time}} \cdot m_{ic}^{\text{Freq}}$  (multiplication term) 1022 and  $m_{ic}^{\text{Time}} \cdot m_{ic}^{\text{Freq}} + m_{i\Omega}^{\text{Time}} \cdot m_{ic}^{\text{Freq}}$  (addition term); (2) Dempster's rule considers additionally the 1023 uncertainty  $m_{i\Omega}^{\text{Time}}$  and  $m_{i\Omega}^{\text{Freq}}$  (in term  $m_{i\Omega}^{\text{Time}} \cdot m_{ic}^{\text{Freq}} + m_{i\Omega}^{\text{Time}} \cdot m_{i\Omega}^{\text{Freq}}$ ), when calculating  $m_{ic}^{\text{Comb}}$ . 1024 To further illustrate the contribution of combining the results from time and frequency domains, we 1025 show the cosine distance between  $\mathbf{g}_{i}^{\text{time}}$  and  $\mathbf{g}_{i}^{\text{Freq}}$  of the same sample in Fig.4(a). "W/o Comb" 1026 means that the cosine distance is calculated between the frequency embeddings learned from "W/o 1026  $\mathcal{L}^{\text{Time}}$  and the time embeddings learned from "W/o  $\mathcal{L}^{\text{Freq}}$ ". One can find that the cosine distance 1027 is smaller by combining the results from time and frequency domains. It means that combining the 1028 results from these two domains indeed enforces the time- and frequency-based embeddings closer 1029 to each other. Taking the FD-A dataset as an example, the average cosine distance decreases from 1.924 to 1.557, i.e., 19.1% by considering the combination.

1031 Compared with TIC-cau and TIC-bol, TIC performs better on all the benchmark datasets. The rea-1032 sons are: the cautious conjunctive rule uses the intersection operation and loses useful information; 1033 the bold conjunctive rule uses the union operation and takes into account confusing information 1034 when making a decision. Comparing TIC with TIC-tri, TIC-tri only has the same ARI as TIC in 2 1035 cases but has lower ARI on the other 7 cases. It is because that the parametric triangular-norm-based 1036 rule is more sensitive to the hyper-parameters, and reasonable hyper-parameter values are difficult to choose in both of the time and frequency domains. This also suggests that the mass functions 1037 associated with the clustering results from the time and frequency domains are independent of each 1038 other in the time series clustering problem studied in our paper. Besides, only TIC-comc in the 5 1039 variants aiming to tackle high-conflict have higher ARI than TIC on SleepEEG dataset, because the 1040 conflict values between the results of time and frequency domains are small. 1041

- 1042 Other 5 variants
- TIC-SVM, TIC-LR, TIC-FC and TIC-LSTM: the outputs from time- and frequency- domains are treated as features to train SVM, Logistic Regression, fully connected layers, and LSTM;
- TIC-max: the maximum probability given in networks from time- and frequency- domains is chosen as the final output probability.

are considered to show that using evidence theory to combine the results from time and frequency domains are superior to other fusion methods. As shown in Table 7, Using evidence theory has the best ARI in 7 of 9 cases, showing that fusing the results via evidence theory is better than other methods.

Table 7: ARI values (mean $\pm$ std.deviation) of different variants of TIC.

ARI	EMG	ECG	HAR	Gesture	FD-A	FD-B	SleepEEG	Epilepsy	UCR	Average
TIC-SVM	.987	.875	.779	.930	.928	.842	.882	.962	.869	
TIC-LR	1	.839	.782	.929	.928	.844	.885	.974	.870	
TIC-FC	1	.860	.769	.934	.932	.859	.876	.980	.863	
TIC-LSTM	.993	.853	.739	.945	.917	.851	.879	.974	.876	
TIC-max	.976	.842	.754	.928	.921	.809	.882	.956	.879	
TIC (Full model)	1	<u>.879</u>	<u>.796</u>	.931	<u>.939</u>	.853	<u>.896</u>	<u>.980</u>	.883	

1060 1061 1062

1053

## A.8 EMBEDDING EVALUATION: TIME-SERIES ANOMALY DETECTION.

1063 We evaluate how TIC performs on a sample-level anomaly detection task, which aims to detect abnormal time-series samples. We build two subsets of FD-B and ECG datasets. The former contains 1064 1000 samples where 900 undamaged bearings are considered "normal" and 100 damaged samples are "outliers"; while the latter has 2000 samples where 1800 normal sinus rhythm recordings are 1066 considered "normal" and 200 atrial fibrillation recordings are "outliers". We randomly select half 1067 of the "normal" samples as the training data, of which the embeddings learned by different models 1068 are used to train the one-class SVM. The time- and frequency-based embeddings learned by TIC are 1069 also concatenated  $[\mathbf{g}^{\text{Time}}; \mathbf{g}^{\text{Freq}}]$ . In Table 8, we report the performance on the anomaly detection 1070 task in terms of Precision, Recall, F1-score and AUROC. TIC achieves the best performance in 6 1071 out of 8 cases. On the ECG dataset, TIC outperforms the second-best model (i.e., activity2vec) 1072 by a margin of 3.3% in AUROC. It shows that TIC can effectively detect the abnormal samples in 1073 mechanical devices and ECGs.

1074

# 1075 A.9 EFFECT OF AUGMENTATION 1076

As shown in Section 3, data augmentations are adopted in both the time- and frequency-based contrastive modules. We explore the effects of augmentation techniques in these two modules separately. In each trial, the corresponding setting is changed while the other ones are the same as the default settings.

Table 8: Performance on time-series anomaly detection. The subsets of FD-B (1000 samples) and ECG (2000 samples) are used. These two subsets are highly imbalanced (90% normal samples and 10% abnormal samples).

FD-B	Precision	Recall	F1-score	AUROC
activity2vec	$.8054 \pm .034$	$.7145 \pm .037$	$.7432 \pm .023$	$.7911_{\pm.029}$
BTSF	$.6854 \pm .015$	$.6274 \pm .017$	$.6653 \pm .021$	$.6741 \pm .014$
CLOCS	$.8412 \pm .025$	$.7465 \pm .028$	$.8222 \pm .019$	.8126 + .021
MHCCL	$.6210 \pm .043$	$.5745 \pm .036$	$.6009 \pm .040$	$.6123 \pm .023$
TFC	<u>.8547</u> ±.037	.7698±.029	<u>.8274</u> ±.019	.8602±.033
Triplet	$.7321 \pm .034$	$.6547_{+.045}$	$.6987_{+.043}$	$.7201 \pm .031$
TS2Vec	$.6813 \pm .022$	$.6134 \pm .027$	$.6423 \pm .019$	$.6731 \pm .015$
TSTCC	$.6354 \pm .019$	$.4517 \pm .066$	$.4968 \pm .056$	$.8022 \pm .044$
TIC (ours)	.8641 ± .025	.7652±.013	.8359±.024	.8563±.033
ECG	Precision	Recall	F1-score	AUROC
activity2vec	$.7584_{\pm.011}$	$.6124_{\pm.036}$	$.7022 \pm .073$	.7598±.026
BTSF	.7752±.027	.7124±.034	$.7413 \pm .022$	$.7581 \pm .029$
CLOCS	$.4861 \pm .049$	$.4035 \pm .056$	$.4491 \pm .038$	$.4727 \pm .023$
MHCCL	$.5689 \pm .015$	$.4875 \pm .019$	$.5471 \pm .024$	.5563 + .023
TFC	$.7654 \pm .037$	$.7035 \pm .029$	.7503±.019	$.7429 \pm .033$
Triplet	$.5789 \pm .027$	.5067 + .023	$.5417 \pm .028$	$.5561 \pm .019$
TS2Vec	$.6857 \pm .015$	$.6138 \pm .019$	$.6587 \pm .021$	$.6701 \pm .020$
TSTCC	$.7345 \pm .031$	$.6538 \pm .024$	$.6993 \pm .016$	$.7235 \pm .019$
TIC (ours)	7958 001	.7216	.7645	.7856

In the time-based contrastive module, we fix the augmentation as jittering, scaling and time-shift for all samples, instead of randomly choosing one augmentation from the augmentation bank  $\mathcal{B}^{\text{Time}}$ (consisting of these three augmentations) for each sample. The comparison result is shown in the rightmost sub-table of Table 9. As can be seen, random selection covers diverse augmentations that allow the encoder  $H_{\text{T}}(\cdot)$  to learn better and more robust embedding  $\mathbf{g}^{\text{Time}}$ , which improves the clustering performance.



Figure 7: T-SNE visualization of the concatenated embedding  $[\mathbf{g}_i^{\text{Time}}; \mathbf{g}_i^{\text{Freq}}]$  for EMG dataset in different epochs. As the learning proceeds, the embeddings from the same cluster are gradually grouped together.

1116

1084

1093 1094 1095

We consider three types of change in frequency-based augmentations. (1) We explore the number of 1117 components manipulated, where the results are shown by the sub-table titled by *Components* in Table 1118 9. One can find that the performance of TIC model tends to deteriorate when perturbing multiple 1119 frequency components. This degradation is attributed to the substantial changes in the time-domain 1120 of augmented samples. Consequently, these augmented samples become readily distinguishable by a 1121 contrastive module, resulting in suboptimal contrastive encoders. (2) We then explore the adjustment 1122 size of amplitude in sub-table Amplitude of Table 9. As can be seen, manipulating the amplitude does 1123 not significantly affect the performance of TIC model. (3) In sub-table Bands of Table 9, we explore 1124 the band that the manipulated component belongs to. The Low-/high-frequency band corresponds 1125 to the first/second half of the frequency spectrum, and contributes to slow/fast variations in the 1126 time domain. In a physiological time-series dataset (e.g., SleepEEG), the low band contains most 1127 information and manipulating a low-frequency component leads to higher ARI. On the contrary, high-band components are more informative than low-band ones in a mechanical time-series dataset 1128 (e.g., FD-A). Thus, perturbing high-band components outperform low-band augmentations. 1129

1130

1132

1131 A.10 VISUALIZATION OF CLUSTERING RESULTS

Figure 7 shows the t-SNE Van der Maaten & Hinton (2008) plot of the embeddings for EMG dataset learned by TIC model. For each sample, we concatenate the embeddings  $[\mathbf{g}_i^{\text{Time}}; \mathbf{g}_i^{\text{Freq}}]$ . As the Table 9: ARI with different augmentation techniques. Terms *Components*, *Amplitudes* and *Bands* refer to the frequency-based augmentations. *Components* means how many components are manipulated, *Amplitude* means the adjustment size of amplitude and *Bands* means the perturbation is performed on a low- or a high-frequency component. *Time domain* means the augmentation  $\tilde{\mathbf{x}}_i^{\text{Time}}$ are randomly selected from the bank  $\mathcal{B}_i^{\text{Time}}$  {Jittering, Scaling Shift} or a fixed augmentation technique.

-	ARI	Components			Amplitude				Bands		Time domain			
		$\beta = 1$	$\beta = 3$	$\beta = 5$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 0.9$	$\gamma = 1.1$	Low	High	Random	Jittering	Scaling	Shift
	SleepEEG	0.8961	0.8751	0.8245	0.8952	<u>0.8961</u>	0.8934	0.8942	<u>0.8994</u>	0.8802	<u>0.8961</u>	0.8726	0.8542	0.8793
_	FD-A	<u>0.9392</u>	0.9157	0.9013	0.9406	0.9392	<u>0.9410</u>	0.9375	0.9321	<u>0.9457</u>	<u>0.9392</u>	0.9103	0.9261	0.9245

learning proceeds, the embeddings learned by TIC from the same cluster are gradually grouped
together. In particular, samples clearly form 3 clusters at epoch = 99, corresponding to the ARI=1
of TIC shown in Table 5. The visualization results further demonstrate the well embedding ability
of our model.