# A GRADUAL COARSE-TO-FINE FRAMEWORK FOR IR REGULARLY SAMPLED MULTIVARIATE TIME SERIES ANALYSIS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

Paper under double-blind review

### ABSTRACT

Irregularly sampled multivariate time series (ISMTS) are prevalent in reality. Most existing methods treat ISMTS as synchronized regularly sampled time series with missing values, neglecting that the irregularities are primarily attributed to variations in sampling rates. In this paper, we introduce a novel perspective that irregularity is essentially relative in some senses. With sampling rates artificially determined from low to high, an irregularly sampled time series can be transformed into a hierarchical set of relatively regular time series from coarse to fine. We observe that additional coarse-grained relatively regular series not only mitigate the irregularly sampled challenges but also incorporate broad-view temporal information, thereby serving as a valuable asset for representation learning. Therefore, following the philosophy of learning that Seeing the big picture first, then delving into the details, we present the Multi-Scale and Multi-Correlation Attention Network (MuSiCNet) combining multiple scales to iteratively refine the ISMTS representation. Specifically, within each scale, we explore time attention and frequency correlation matrices to aggregate intra- and inter-series information, naturally enhancing the representation quality with richer and more intrinsic details. While across adjacent scales, we employ a representation rectification method containing contrastive learning and reconstruction results adjustment to further improve representation consistency. MuSiCNet is an ISMTS analysis framework that competitive with SOTA in three mainstream tasks consistently, including classification, interpolation, and forecasting.

1 INTRODUCTION

Irregularly sampled multivariate time series (ISMTS) are ubiquitous in realistic scenarios, ranging from scientific explorations to societal interactions (Che et al., 2018; Shukla & Marlin, 2021; Sun 037 et al., 2021; Agarwal et al., 2023; Yalavarthi et al., 2024). The causes of irregularities in time series collection are diverse, including sensor malfunctions, transmission distortions, cost-reduction strategies, and various external forces or interventions, etc. Such ISMTS data exhibit distinctive 040 features including intra-series irregularity, characterized by inconsistent intervals between consecutive 041 data points, and inter-series irregularity, marked by a lack of synchronization across multiple variables. 042 The above characteristics typically result in the lack of alignment and uneven count of observations 043 (Shukla & Marlin, 2020), invalidating the assumption of coherent fixed-dimensional feature space for 044 most traditional time series analysis models.

Recent studies have attempted to address these challenges by treating ISMTS as synchronized, regularly sampled Normal Multivariate Time Series (NMTS) data with missing values, focusing on imputation strategies (Che et al., 2018; Yoon et al., 2018; Camino et al., 2019; Tashiro et al., 2021; Zhang et al., 2021c; Chen et al., 2022; Fan, 2022; Du et al., 2023). However, direct imputation is difficult, especially when sampling is sparse. Inaccurate imputation results can distort underlying relationships and introduce significant noise, which can greatly reduce the accuracy of analysis tasks (Zhang et al., 2021b; Wu et al., 2021; Agarwal et al., 2023; Sun et al., 2024). Latest developments circumvent imputation and aim to address these challenges by embracing the inherent continuity of time, thus preserving the continuous temporal dynamics dependencies within the ISMTS data. Despite these innovations, most methods above are merely solutions for intra-series irregularities,

- such as Recurrent Neural Networks (RNNs) (De Brouwer et al., 2019; Schirmer et al., 2022; Agarwal
  et al., 2023)- and Neural Ordinary Differential Equations (Neural ODEs)-based methods (Kidger
  et al., 2020; Rubanova et al., 2019; Jhin et al., 2022; Jin et al., 2022) and the unaligned challenges
  presented by inter-series irregularities in multivariate time series remain unsolved.
- 058 Delving into the nature of irregularly sampled time series, we discover that the intra- and inter-060 series irregularities in ISMTS primarily arise 061 from inconsistency in sampling rates within and 062 across variables. We argue that irregularities are 063 essentially relative in some senses and by arti-064 ficially determined sampling rates from low to high, ISMTS can be transformed into a hierarchi-065 cal set of relatively regular time series from coarse 066 to fine. Taking a broader perspective, setting a 067 lower and consistent sampling rate within an in-068 stance can synchronize sampling times across se-069 ries and establish uniform time intervals within series. This approach can mitigate both types 071 of irregularity and emphasize long-term dependencies. As shown in Fig.1, the coarse-grained 073 scales 1 and 2 exhibit balanced placements for 074 all variables in the instance and provide clearer 075 overall trends. However, lower sampling rates may lead to information loss and sacrifice de-076 tailed temporal variations. Conversely, with a 077 higher sampling rate as in scale L, more real observations contain rich information and prevent 079 artificially introduced dependencies beyond original relations during training. Nonetheless, the 081

101

102

103

104

105

106



Figure 1: Comparative visualization of multiscale time series data with various sampling rates. Scale L depicts the original selected representative time series in the P12 Dataset to show the inter- and intra-series irregularities. Scale 1 to Scale L - 1 illustrates the effect of applying different sampling rates from low to high.

- significant irregularity in fine-grained scales poses a greater challenge for representation learning.
- To bridge this gap, we propose MuSiCNet-a Multi-Scale and Multi-Correlation Attention Net-083 work—to iteratively optimize ISMTS representations from coarse to fine. Our approach begins by 084 establishing a hierarchical set of coarse- to fine-grained series with sampling rates from low to high. 085 At each scale, we employ a custom-designed encoder-decoder framework called multi-correlation attention network (CorrNet), for representation learning. Different from most existing methods that 087 focus mainly on intra-series relationships, our CorrNet encoder (CorrE) captures embeddings of 880 continuous time values by using an attention mechanism containing correlation matrices to aggregate both intra- and inter-series information. This approach is crucial not only because every observation 090 in ISMTS, given the sparse sampling, is valuable for representation learning, but also due to the fact 091 that correlated variables provide deeper insights for a given query. Therefore, we designed frequency 092 correlation matrices using Lomb-Scargle Periodogram-based Dynamic Time Warping (LSP-DTW). This approach addresses the challenges of calculating correlations in ISMTS and re-weights the inter-series attention scores to better capture cross-series information. Across scales, we employ a 094 representation rectification operation from coarse to fine to iteratively refine the learned representa-095 tions with contrastive learning and reconstruction results adjustment methods. This ensures accurate 096 and consistent representation and minimizes error propagation throughout the model.
- Benefiting from the aforementioned designs, MuSiCNet explicitly learns multi-scale information,
   enabling good performance on widely used ISMTS datasets, thereby demonstrating its ability to
   capture relevant features for ISMTS analysis. Our main contributions can be summarized as follows:
  - We find that irregularities in ISMTS are essentially relative in some senses and multi-scale learning helps balance coarse- and fine-grained information in ISMTS representation learning.
  - We introduce CorrNet, an encoder-decoder framework designed to learn fixed-length representations for ISMTS. Notably, our proposed LSP-DTW can mitigate spurious correlations induced by irregularities in the frequency domain and effectively re-weight attention across sequences.
- We are not limited to a specific analysis task and attempt to propose a task-general model for ISMTS analysis, including classification, interpolation, and forecasting.



Figure 2: Overview of MuSiCNet framework, shown in (a), containing three main components for better representation learning, including hierarchical structure  $\{X_{\text{mask}}^{(l)}\}_{l=1}^{L}$ , representation learning using CorrNet within scale  $\ell_{\text{cons}}^{(l)}$ , and rectification operation across adjacent scales  $\ell_{\text{recon}}^{(l)}$ . (b) visualizes the encoding process in CorrNet for Scale *l*, which relies on  $\tau_n^{(l)}$  to aggregates intra-series information, and then relies on  $c_{d_i,(\cdot)}$  to fuse inter-series information from other variables for  $d_i$ -th dimension. (c) visualizes the calculation process of the correlation matrix, which transfers the time domain into the frequency domain with LSP, and then utilizes DTW to calculate the similarity weight.

136 137

138

## 2 RELATED WORK

139 Irregularly Sampled Multivariate Time Series Analysis. An effective approach for analyzing ISMTS hinges on the understanding of their unique properties. Most existing methods treat ISMTS 140 as NMTS with missing values, such as Che et al. (2018); Yoon et al. (2018); Camino et al. (2019); 141 Tashiro et al. (2021); Chen et al. (2022); Fan (2022); Du et al. (2023); Wang et al. (2024). However, 142 most imputation-based methods may distort the underlying relationships, introducing unsuitable 143 inductive biases and substantial noise due to incorrect imputation (Zhang et al., 2021b; Wu et al., 144 2021; Agarwal et al., 2023), ultimately compromising the accuracy of downstream tasks. Some 145 other methods treat ISMTS as time series with discrete timestamps, aggregating all sample points 146 of a single variable to extract a unified feature for each variable (Zhang et al., 2021b; Horn et al., 147 2020; Li et al., 2023). These methods can directly accept raw ISMTS data as input but often struggle 148 to handle the underlying relationships within the time series. Recent progress seeks to overcome 149 these challenges by recognizing and utilizing the inherent continuity of time, thereby maintaining the ongoing temporal dynamics present in ISMTS data (De Brouwer et al., 2019; Rubanova et al., 2019; 150 Kidger et al., 2020; Schirmer et al., 2022; Jhin et al., 2022; Chowdhury et al., 2023). 151

Despite these advancements, existing methods mainly suffer from two main drawbacks, they primarily address intra-series irregularity while overlooking the alignment issues stemming from inter-series irregularity, and 2) they rely on assumptions tailored to specific downstream tasks (Yalavarthi et al., 2024; Wang et al., 2024), hindering their ability to consistently perform well across various ISMTS tasks.

Multi-scale Modeling. Multi-scale and hierarchical approaches have demonstrated their utility across various fields, including computer vision (CV) (Fan et al., 2021; Zhang et al., 2021a), natural language processing (NLP) (Nawrot et al., 2021; Zhao et al., 2021), and time series analysis (Chen et al., 2021; Shabani et al., 2022; Cai et al., 2024). Most recent innovations in the time series analysis domain have seen the integration of multi-scale modules into the Transformer architecture to enhance analysis capabilities Shabani et al. (2022); Liu et al. (2021) and are designed for NMTS.

162 Nevertheless, the application of multi-scale modeling specifically designed for ISMTS data, and 163 the exploitation of information across scales, remain less unexplored. As far as we know, Singh 164 et al. (2019) and Zhang et al. (2023) are among the earlier works on multi-level ISMTS learning. 165 Singh et al. (2019) addresses multi-resolution signal issues by distributing signals across specialized 166 branches with different resolutions, where each branch employs a Flexible Irregular Time Series Network (FIT) to process high- and low-frequency data separately. Zhang et al. (2023), on the other 167 hand, is a transformer-based model that stacks multiple Warpformer layers to produce multi-scale 168 representations, combining them via residual connections to support downstream tasks. These works typically focus on either specific tasks or particular model architectures. In contrast, our design 170 philosophy originates from ISMTS characteristics rather than being tied to a specific feature extraction 171 network structure. Warpformer emphasizes designing a specific network architecture but involves 172 high computational costs and requires manually balancing the trade-off between the number of scales 173 and the dataset. These are challenges that our MuSiCNet avoids entirely.

174 175 176

177

## 3 PROPOSED MUSICNET FRAMEWORK

As previously mentioned, our work aims to learn ISMTS representation for further analysis tasks by introducing MuSiCNet, a novel framework designed to balance coarse- and fine-grained information across different scales. The overall model architecture illustrated in Fig.2(a) indicates the effectiveness of MuSiCNet can be guaranteed to a great extent by 1) Hierarchical Structure. 2) Representation Learning Using CorrNet Within Scale. 3) Rectification Across Adjacent Scales. We will first introduce problem formulation and notations of MuSiCNet and then discuss key points in the following subsections.

184 185

186

## 3.1 PROBLEM FORMULATION

187 Our goal is to learn a nonlinear embedding function  $f_{\theta}$ , such that the set of ISMTS data  $\mathcal{X} = \{X_1, \dots, X_N\}$  can map to the best-described representations for further ISMTS analysis including 189 both supervised and unsupervised tasks. We denote  $X_n \in \mathbb{R}^{T_n \times D}$  as a D-dimensional instance with 190 the length of observation  $T_n$ . Specifically, the *d*-th dimension in instance *n* can be treated as a tuple 191  $X_{dn} = (x_{dn}, t_{dn})$  where the length of observations is  $T_{dn}$ .  $x_{dn} = [x_{1dn}, \dots x_{T_{dn}dn}]$  is the list of 192 observations and the list of corresponding observed timestamps is  $t_{dn} = [t_{1dn}, \dots t_{T_{dn}dn}]$ . We drop 193 the data case index *n* for brevity when the context is clear.

194 195

## 3.2 CORRNET ARCHITECTURE WITHIN SCALE

Multi-Correlation Attention Module. In this subsection, we elaborate on the Multi-Correlation Attention module. Time attention has proven effective for ISMTS learning (Shukla & Marlin, 2021; Horn et al., 2020; Chowdhury et al., 2023; Yu et al., 2024). Many existing methods mainly focus on capturing interactions between observation values and their corresponding sampling times within a single variable. However, due to the potential sparse sampling in ISMTS, observations from all variables are valuable and need to be considered.

To address this, we use irregularly sampled time points and corresponding observations from all variables within a sample as keys and values to produce fixed-dimensional representations at the query time points. The importance of each variable cannot be uniform for a given query and variables that provide more valuable information should receive more attention. Therefore, we designed frequency correlation matrices to re-weight the inter-series attention scores, enhancing the representation learning process.

In general, as illustrated in Fig.2(b), taking ISMTS X as input, the CorrNet Encoder CorrE (·) generates multi-time attention embedding as follows:

$$CorrE(\boldsymbol{Q}_T, \boldsymbol{K}_T, \boldsymbol{X}) = \boldsymbol{A}_T \boldsymbol{X} \boldsymbol{C}_T$$
  
$$\boldsymbol{A}_T = \text{softmax}(\boldsymbol{Q}_T \boldsymbol{K}_T / d_r)$$
(1)

211 212 213

where the calculation of  $A_T \in \mathbb{R}^{K \times T}$  is based on a time attention mechanism with query  $Q_T \in \mathbb{R}^{K \times K}$  and key  $K_T \in \mathbb{R}^{K \times T}$  (Vaswani et al., 2017). Since more attention should be paid to correlated variables for a given query which can provide more valuable knowledge. Therefore,

different input dimensions should utilize various weights of time embeddings through the correlation matrix  $C_T \in \mathbb{R}^{D \times D}$ , and we will introduce it in the following Correlation Extraction paragraph.

Since the continuous function defined by the CorrE module is incompatible with neural network architectures designed for fixed-dimensional vectors or discrete sequences, following the method in Shukla & Marlin (2021), we generate an output representation by materializing its output at a predefined set of reference time points  $\tau = [\tau_1, \dots, \tau_K]$ . This process transforms the continuous output into a fixed-dimensional vector or a discrete sequence, thereby making it suitable for subsequent neural network processing.

224 225

**Correlation Extraction.** The correlation matrix is essential for deriving reliable and consistent 226 correlations within ISMTS, which must be robust to the inherent challenges of variable sampling 227 rates and inconsistent observation counts at each timestamp in ISMTS. Most existing distance 228 measures, such as Euclidean distance, Dynamic Time Warping (DTW) (Berndt & Clifford, 1994), and 229 Optimal Transport / Wasserstein Distance (Villani et al., 2009), risk generating spurious correlations 230 in the context of irregularly sampled time series. This is due to their dependence on the presence 231 of both data points for the similarity measurement, and the potential for imputation to introduce 232 unreliable information before calculating similarity which will be explored further in Section 4.4 of 233 our experiments.

234 At an impasse, the Lomb-Scargle Periodogram (LSP) (Lomb, 1976; Scargle, 1982) provides en-235 lightenment to address this issue. LSP is a well-known algorithm for generating a power spectrum 236 and detecting the periodic component in irregularly sampled time series. It extends the *Fourier* 237 periodogram approach to accommodate irregularly sampled scenarios (VanderPlas, 2018) eliminat-238 ing the need for interpolation or imputation. This makes LSP a great tool for simplifying ISMTS 239 analysis. Compared to existing methods, measuring the similarity between discrete raw observations, LSP-DTW, an implicit continuous method, utilizes inherent periodic characteristics and provides 240 global information to measure the similarity. 241

As demonstrated in Fig.2(c), we first convert ISMTS into the frequency domain using LSP and then apply DTW to evaluate the distance between variables. The correlation between  $X_{d_i}$  and  $X_{d_j}$  is:

$$c_{d_i d_j} = \mathsf{DTW}\left(\mathsf{LSP}(\boldsymbol{X}_{d_i}), \mathsf{LSP}(\boldsymbol{X}_{d_j})\right) = \min_{\pi} \sum_{(m,n) \in \pi} \left(\mathsf{LSP}(\boldsymbol{X}_{d_i})[m] - \mathsf{LSP}(\boldsymbol{X}_{d_j})[n]\right)^2$$
(2)

where  $\pi$  is the search path of DTW. We calculate the correlation matrix  $C_T$  by iteratively performing the aforementioned step for an instance.

Notably, we compute the correlation matrix using LSP-DTW only once per instance, without iteratively applying it, and it is not calculated in model training or inference.

252

245 246

253 **Encoder-Decoder Framework.** Drawing inspiration from notable advances in NLP and CV, our 254 core network, CorrNet employs time series masked modeling, which learns effective time series 255 representations to facilitate various downstream analysis tasks. It is a framework consisting of an encoder-decoder architecture based on continuous-time interpolation. At each scale l, CorrE 256 learns a set of latent representations  $r^{(l)} = [r_1, \cdots, r_K]$  defined at K reference time points on the 257 randomly masked ISMTS. We further employ CorrNet Decoder (CorrD), a simplified CorrE (without 258 correlation matrix), to produce the reconstructed output  $\hat{X}_{reco}^{(l)}$ , using the input time point sequence 259 260  $t^{(l)}$  as reference points. We iteratively apply the same CorrNet at each scale. Here, we emphasize that all scales share a single encoder that can reduce the model complexity and keep feature extraction 261 consistency for various scales. 262

We measure the reconstruction accuracy using the Mean Squared Error (MSE) between the reconstructed values and the original ones at each timestamp and calculate the MSE loss specifically for the masked timestamps, as expressed in the following equation

$$\ell_{\text{recon}}^{(l)} = \sum_{n=1}^{N} \left\| \boldsymbol{M}^{(l)} \odot \left( \left( \hat{\boldsymbol{X}}_{\text{reco}}^{(l)} \right)_{n} - \boldsymbol{X}_{n}^{(l)} \right) \right\|_{2}^{2}$$
(3)

267 268 269

where  $M^{(l)}$  is the mask for the *l*-th scale,  $\odot$  is the Hadamard product.

# 270 3.3 RECTIFICATION STRATEGY ACROSS SCALES271 3.3 RECTIFICATION STRATEGY ACROSS SCALES

Following the principle that adjacent scales exhibit similar representations and coarse-grained scales
contain more long-term information, the rectification strategy is a key component of our MuSiCNet
framework. We implement a dual rectification strategy across adjacent scales to enhance representation learning. We start by generating a hierarchical set of relatively regular time series from coarse to
fine by

$$\boldsymbol{X}_{\text{multi}} = \boldsymbol{M}^{(l)} \odot (\text{AvgPooling}_{L}(\boldsymbol{X})) = \{\boldsymbol{X}_{\text{mask}}^{(1)}, \cdots, \boldsymbol{X}_{\text{mask}}^{(L)}\}$$
(4)

While the coarse-grained series ignores detailed variations for high-frequency signals and focuses on much clearer broad-view temporal information, the fine-grained series retains detailed variations for frequently sampled series. As a result, iteratively using coarse-grained information for fine-grained series as a strong structural prior can benefit ISMTS learning.

Firstly, the reconstruction results at scale l is designed to align closely with the results at the (l-1)-th scale, that is to say, the reconstruction results at scale (l-1) can be used to adjust the results at scale l using MSE,

$$\ell_{\text{adj}}^{(l)} = \sum_{n=1}^{N} \left\| \left( \text{AvgPooling}_{l}(\hat{\boldsymbol{X}}_{\text{reco}}^{(l)}) \right)_{n} - \left( \hat{\boldsymbol{X}}_{\text{reco}}^{(l-1)} \right)_{n} \right\|_{2}^{2}$$
(5)

Secondly, contrastive learning is leveraged to ensure coherence between adjacent scales. Pulling these two representations between adjacent scales together and pushing other representations within the batch  $\mathcal{B}$  apart, not only facilitates the learning of within-scale representations but also enhances the consistency of cross-scale representations. Taking into consideration that the dimensions of  $r^{(l)}$  and  $r^{(l-1)}$  are different, we employ a GRU Network as a decoder to uniform dimension as  $h^{(l)}$  and  $h^{(l-1)}$  before contrastive learning.

293 294 295

296

298

299

300

301

302

286 287

288

289

290

291

292

277

$$\ell_{\text{cons}}^{(l)} = -\sum_{i=1}^{N} \log \frac{\exp\left(\boldsymbol{h}_{i}^{(l)} \cdot \boldsymbol{h}_{i}^{(l-1)}\right)}{\sum_{j=1}^{\mathcal{B}} \left(\exp\left(\boldsymbol{h}_{i}^{(l)} \cdot \boldsymbol{h}_{j}^{(l-1)}\right) + \mathbb{I}_{[i \neq j]} \exp\left(\boldsymbol{h}_{i}^{(l)} \cdot \boldsymbol{h}_{j}^{(l)}\right)\right)}$$
(6)

297 where the  $\mathbb{I}$  is the indicator function.

The advantage of the two operations lies in their ability to ensure a consistent and accurate representation of the data at different scales. This strategy significantly improves the model's ability to learn representations from ISMTS data, which is essential for tasks requiring detailed and accurate time series analysis. Last but not least, this method ensures that the model remains robust and effective even when dealing with data at varying scales, making it versatile for diverse applications.

303 304 305

306

307

308

309

310

## 4 EXPERIMENT

In this section, we demonstrate the effectiveness of the MuSiCNet framework for time series classification, interpolation and forecasting. *Notably, for each dataset, the window size is initially set to* 1/4 *of the time series length and then halved iteratively until the majority of the windows contain at least one observation.* Our results are based on the mean and standard deviation values computed over 5 independent runs. **Bold** indicates the best performer, while <u>underline</u> represents the second best. Due to the page limitation, we provide more detailed setup for experiments in the Appendix.

311 312 313

314

## 4.1 TIME SERIES CLASSIFICATION

Datasets and experimental settings. We use real-world datasets including healthcare and human activity for classification. (1) P19 (Reyna et al., 2020) with missing ratio up to 94.9%, includes 38, 803 patients that are monitored by 34 sensors. (2) P12 (Goldberger et al., 2000) records temporal measurements of 36 sensors of 11, 988 patients in the first 48-hour stay in ICU, with a missing ratio of 88.4%. (3) PAM (Reiss & Stricker, 2012) contains 5, 333 segments from 8 activities of daily living that are measured by 17 sensors and the missing ratio is 60.0%. *Importantly, P19 and P12 are imbalanced binary label datasets*.

Here, we follow the common setup by randomly splitting the dataset into training (80%), validation (10%), and test (10%) sets and the indices of these splits are fixed across all methods. Consistent with prior researches, we evaluate the performance of our framework on classification tasks using the area

326		P19		P	P12		PAM				
327	Methods	AUROC	AUPRC	AUROC	AUPRC	Accuracy	Precision	Recall	F1 score		
328	Transformer	$80.7 \pm 3.8$	$42.7 \pm 7.7$	$83.3 \pm 0.7$	$47.9 \pm 3.6$	$83.5 \pm 1.5$	$84.8 \pm 1.5$	$86.0 \pm 1.2$	$85.0 \pm 1.3$		
329	Trans-mean	$83.7 \pm 1.8$	$45.8 \pm 3.2$	$82.6 \pm 2.0$	$46.3 \pm 4.0$	$83.7 \pm 2.3$	$84.9 \pm 2.6$	$86.4 \pm 2.1$	$85.1 \pm 2.4$		
000	GRU-D	$83.9 \pm 1.7$	$46.9 \pm {\scriptstyle 2.1}$	$81.9 \pm 2.1$	$46.1 \pm 4.7$	$83.3 \pm 1.6$	$84.6 \pm 1.2$	$85.2 \pm 1.6$	$84.8 \pm 1.2$		
330	SeFT	$81.2 \pm 2.3$	$41.9 \pm 3.1$	$73.9 \pm 2.5$	$31.1 \pm 4.1$	$67.1 \pm 2.2$	$70.0 \pm 2.4$	$68.2 \pm 1.5$	$68.5 \pm 1.8$		
331	mTAND	$84.4 \pm 1.3$	$50.6 \pm 2.0$	$84.2 \pm 0.8$	$48.2 \pm 3.4$	$74.6 \pm 4.3$	$74.3 \pm 4.0$	$79.5 \pm 2.8$	$76.8 \pm 3.4$		
	IP-Net	$84.6 \pm 1.3$	$38.1 \pm 3.7$	$82.6 \pm 1.4$	$47.6 \pm 3.1$	$74.3 \pm 3.8$	$75.6 \pm 2.1$	$77.9 \pm 2.2$	$76.6 \pm 2.8$		
332	DGM <sup>2</sup> -O	$86.7 \pm 3.4$	$44.7 \pm 11.7$	$84.4 \pm 1.6$	$47.3 \pm 3.6$	$82.4 \pm 2.3$	$85.2 \pm 1.2$	$83.9 \pm 2.3$	$84.3 \pm 1.8$		
333	MTGNN	$81.9 \pm 6.2$	$39.9 \pm 8.9$	$74.4 \pm 6.7$	$35.5 \pm 6.0$	$83.4 \pm 1.9$	$85.2 \pm 1.7$	$86.1 \pm 1.9$	$85.9 \pm 2.4$		
	Raindrop	$87.0 \pm 2.3$	$51.8 \pm 5.5$	$82.8 \pm 1.7$	$44.0 \pm 3.0$	$88.5 \pm 1.5$	$89.9 \pm 1.5$	$89.9 \pm 0.6$	$89.8 \pm 1.0$		
334	Warpformer	$88.8 \pm 1.7$	$55.2 \pm 3.9$	$83.4 \pm 0.9$	$47.2 \pm 3.7$	$94.3 \pm 0.6$	$95.8 \pm 0.8$	$94.8 \pm 1.0$	$95.2 \pm 0.6$		
335	ViTST	$89.2 \pm 2.0$	$53.1 \pm 3.4$	$\underline{85.1} \pm 0.8$	$\underline{51.1} \pm 4.1$	$\underline{95.8} \pm 1.3$	$\underline{96.2} \pm 1.3$	$\underline{96.1} \pm 1.1$	$\underline{96.5} \pm 1.2$		
336	MuSiCNet	$86.8 \pm 1.4$	$45.4 \pm 2.7$	$86.1 \pm 0.4$	$54.1 \pm 2.2$	$\textbf{96.3} \pm 0.7$	$96.9 \pm 0.6$	$96.9 \pm 0.5$	$96.8 \pm 0.5$		

Table 2: Comparison with the baseline methods on ISMTS interpolation task on PhysioNet.

Model	Mean Squared Error $(\times 10^{-3})$								
Observed %	50%	60%	70%	80%	90%				
RNN-VAE L-ODE-RNN L-ODE-ODE mTAND-Full	$\begin{array}{c} 13.418 \pm 0.008 \\ 8.132 \pm 0.020 \\ 6.721 \pm 0.109 \\ \underline{4.139} \pm 0.029 \end{array}$	$\begin{array}{c} 12.594 \pm 0.004 \\ 8.140 \pm 0.018 \\ 6.816 \pm 0.045 \\ \underline{4.018} \pm 0.048 \end{array}$	$\begin{array}{c} 11.887 \pm 0.005 \\ 8.171 \pm 0.030 \\ 6.798 \pm 0.143 \\ \underline{4.157} \pm 0.053 \end{array}$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{c} 11.470 \pm 0.006 \\ 8.402 \pm 0.022 \\ 7.142 \pm 0.066 \\ \underline{4.798} \pm 0.036 \end{array} $				
MuSiCNet	$0.918 \pm 0.025$	$0.919 \pm 0.064$	$0.938 \pm 0.014$	$0.992 \pm 0.008$	$0.965 \pm 0.008$				

under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPRC) for the P12 and P19 datasets, given their imbalanced nature. For the nearly balanced PAM dataset, we employ Accuracy, Precision, Recall, and F1 Score. For all of the above metrics, higher results indicate better performance.

**Main Results of classification.** We compare MuSiCNet with ten state-of-the-art irregularly sampled time series classification methods, including Transformer (Vaswani et al., 2017), Trans-mean, GRU-D (Che et al., 2018), SeFT (Horn et al., 2020), and mTAND (Shukla & Marlin, 2021), IP-Net (Shukla & Marlin, 2018), DGM<sup>2</sup>-O(Wu et al., 2021), MTGNN (Wu et al., 2020), Raindrop (Zhang et al., 2021b), ViTST (Li et al., 2023) and Warpformer (Zhang et al., 2023). Since mTAND is proven superior over various recurrent models, such as RNNImpute (Che et al., 2018), Phased-LSTM (Neil et al., 2016) and ODE-based models like LATENT-ODE and ODE-RNN (Chen et al., 2018), we focus our comparisons on mTAND and do not include results for the latter model. 

As indicated in Table 1, MuSiCNet demonstrates good performance across three benchmark datasets, underscoring its effectiveness in typical time series classification tasks. Notably, in binary classifica-tion scenarios, MuSiCNet surpasses the best-performing baselines on the P12 dataset by an average of 1.0% in AUROC and 3.0% in AUPRC. For the P19 dataset, while our performance is competitive, MuSiCNet stands out due to its lower time and space complexity compared to ViTST. ViTST converts 1D time series into 2D images, potentially leading to significant space inefficiencies due to the introduction of extensive blank areas, especially problematic in ISMTS. In the more complex task of 8-class classification on the PAM dataset, MuSiCNet surpasses current methodologies, achieving a 0.5% improvement in accuracy and a 0.7% increase in precision. 

Notably, the *consistently low standard deviation* in our results indicates that MuSiCNet is a reliable
 model. Its performance remains steady across varying data samples and initial conditions, suggesting
 a strong potential for generalizing well to new, unseen data. This stability and predictability in
 performance enhance the confidence in the model's predictions, which is particularly crucial in
 sensitive areas such as medical diagnosis in clinical settings.

Methods	USHCN	MIMIC-III	Physionet12
DLinear+	$0.347 \pm 0.065$	$0.691 \pm 0.016$	$0.380 \pm 0.001$
NLinear+	$0.452 \pm 0.101$	$0.726 \pm 0.019$	$0.382 \pm 0.001$
Informer+	$0.320 \pm 0.047$	$0.512 \pm 0.064$	$0.347 \pm 0.001$
FedFormer+	$2.990 \pm 0.476$	$1.100 \pm 0.059$	$0.455 \pm 0.004$
NeuralODE-VAE	$0.960 \pm 0.110$	$0.890 \pm 0.010$	-
GRU-Simple	$0.750 \pm 0.120$	$0.820 \pm 0.050$	-
GRU-D	$0.530 \pm 0.060$	$0.790 \pm 0.060$	-
T-LSTM	$0.590 \pm 0.110$	$0.620 \pm 0.050$	-
mTAND	$0.300 \pm 0.038$	$0.540 \pm 0.036$	$0.315 \pm 0.002$
GRU-ODE-Bayes	$0.430 \pm 0.070$	$0.480 \pm \textbf{0.480}$	$0.329 \pm 0.004$
Neural Flow	$0.414 \pm 0.102$	$0.490 \pm 0.004$	$0.326 \pm 0.004$
CRU	$0.290 \pm 0.060$	$0.592 \pm 0.049$	$0.379 \pm 0.003$
GraFITi	$0.272 \pm 0.047$	$0.396 \pm 0.030$	$0.286 \pm 0.001$
MuSiCNet	$0.268 \pm 0.038$	$\underline{0.475} \pm 0.031$	$0.312 \pm 0.000$

Table 3: Experimental results for **forecasting** next three time steps. – indicates no published results. 

4.2 TIME SERIES INTERPOLATION

**Datasets and experimental settings.** PhysioNet (Silva et al., 2012) consists of 37 variables extracted from the first 48 hours after admission to the ICU. We use all 8,000 instances for interpolation experiments whose missing ratio is 78.0%.

We randomly split the dataset into a training set, encompassing 80% of the instances, and a test set, comprising the remaining 20% of instances. Additionally, 20% of the training data is reserved for validation purposes. The performance evaluation is conducted using MSE, where lower values indicate better performance. 

**Main Results of Interpolation.** For the interpolation task, we compare it with RNN-VAE, L-ODE-RNN (Chen et al., 2018), L-ODE-ODE (Rubanova et al., 2019), mTAND-full. 

For the interpolation task, models are trained to predict or reconstruct values for the entire dataset based on a selected subset of available points. Experiments are conducted with varying observation levels, ranging from 50% to 90% of observed points. During test time, models utilize the observed points to infer values at all time points in each test instance. 

As illustrated in Table 2, MuSiCNet demonstrates superior performance, highlighting its effectiveness in time series interpolation. This can be attributed to its ability to interpolate progressively from coarse to fine, aligning with the intuition of multi-resolution signal approximation (Mallat, 1989).

4.3 TIME SERIES FORECASTING

**Datasets and Experimental Settings.** (1) **USHCN** (Menne et al., 2015) is an artificially prepro-cessing dataset containing measurements of 5 variables from 1280 weather stations in the USA. The missing ratio is 78.0%. (2) MIMIC-III (Johnson et al., 2016) are dataset that rounded the recorded observations into 96 variables, 30-minute intervals and only use observations from the 48 hours after admission. The missing ratio is 94.2%. (3) Physionet12 (Silva et al., 2012) comprises medical records from 12,000 ICU patients. It includes measurements of 37 vital signs recorded during the first 48 hours of admission and the missing ratio is 80.4%. We use MSE to measure forecasting performance, with lower values indicating better performance. 

Main Results of Forecasting. We compare the performance with the ISMTS forecasting models: Graph-based method Grafiti (Yalavarthi et al., 2024), ODE- and RNN-based models including GRU-ODE-Bayes (De Brouwer et al., 2019), Neural Flows (Biloš et al., 2021), CRU (Schirmer et al., 2022), NeuralODE-VAE(Chen et al., 2018), GRUSimple, GRU-D and TLSTM(Baytas et al., 2017). Additionally, attention-based models like mTAND, also an interpolation model, together with variants of Informer (Zhou et al., 2021), Fedformer (Zhou et al., 2022), DLinear, and NLinear (Zeng et al., 2023), denoted as Informer+, Fedformer+, DLinear+, and NLinear+, respectively.



Figure 3: Visualization of various methods to extract the correlation matrix from P12 dataset. The darker the color, the more similar the relationship. (a) denotes the average pairwise observation rate (i.e., 1 minus missing rate), and (b) - (d) denotes different correlation matrices.

Table 4: Classification performance of MuSiCNetr to verify the necessity of correlation matrix.

	P12		P19			PAM		
Methods	AUROC	AUPRC	AUROC	AUPRC	Accuracy	Precision	Recall	F1 score
MuSiCNet w/o Corr Learnable Corr	$ \begin{vmatrix} 86.1 \pm 0.4 \\ 85.5 \pm 0.3 \\ 85.7 \pm 0.4 \end{vmatrix} $	$\begin{array}{c} {\bf 54.1} \pm 2.2 \\ {\bf 53.0} \pm 2.1 \\ {\bf 53.0} \pm 2.0 \end{array}$	$\begin{array}{c c} \textbf{86.8} \pm 0.4 \\ 82.9 \pm 0.8 \\ 83.4 \pm 0.7 \end{array}$	$\begin{array}{c} {\bf 54.1} \pm 2.2 \\ {\bf 32.7} \pm 2.1 \\ {\bf 31.8} \pm 2.7 \end{array}$	$ \begin{vmatrix} 96.3 \pm 0.7 \\ 95.7 \pm 0.9 \\ 96.1 \pm 0.5 \end{vmatrix} $	$\begin{array}{c} \textbf{96.9} \pm 0.6 \\ \textbf{96.2} \pm 0.51 \\ \textbf{96.7} \pm 0.38 \end{array}$	$\begin{array}{c} \textbf{96.9} \pm 0.5 \\ 96.5 \pm 0.2 \\ 96.5 \pm 0.7 \end{array}$	$\begin{array}{c} \textbf{96.8} \pm 0.5 \\ 96.3 \pm 0.3 \\ 96.6 \pm 0.5 \end{array}$

This experiment is conducted following the setting of GraFITi where for the USHCN dataset, the model observes for the first 3 years and forecasts the next 3 time steps and for other datasets, the model observes the first 36 hours in the series and predicts the next 3 time steps.

As shown in Table 3, MuSiCNet consistently achieves competitive performance across all datasets, maintaining accuracy within the top two among baseline models. While GraFITi excels by explicitly modeling the relationship between observation and prediction points, making it superior in certain scenarios, MuSiCNet remains competitive without imposing priors for any specific task.

4.4 CORRELATION RESULTS

In this section, we focus on validating the necessity, effectiveness, and efficiency of the correlation
 matrix in the classification task as an example.

First, we verify the necessity of the correlation matrix using results from all classification datasets in
Table 4. Removing the correlation matrix (line 4) led to performance drops across all datasets, with
P19 showing the largest decline due to its 94.9% missing rate. This highlights the importance of
capturing inter-series relationships in irregularly sampled time series, making the correlation matrix
essential. Replacing the designed correlation matrix with a learnable one (line 5) also worsened
performance, indicating that learning inter-series relationships purely from the network remains
highly challenging and specialized correlation designs are needed.

Second, we evaluate LSP-DTW against other correlation calculation methods (I-GAK (Cuturi, 2011),
I-DTW (Berndt & Clifford, 1994)) on the P12 dataset to verify the effectiveness. Interpolation-based
methods (I-GAK, I-DTW) distort correlations, leading to unreliable results as seen in Fig.3. I-GAK
shows fictitious correlations based on observation rates, while I-DTW presents uniformly positive
correlations, neither of which captures true data characteristics. In contrast, LSP-DTW accurately
identifies correlations, verified by Table 5, where it outperforms all baselines, demonstrating the
importance of appropriate correlation modeling.

Lastly, we report the computation time for the correlation matrix. LSP-DTW based correlation matrix
 is computed per instance in parallel, with acceptable runtimes (0.137s for P12, 0.127s for P19, 0.049s for PAM). It is calculated once, making it efficient for the entire learning process.

486 Table 5: Classification performance of MuSiC-487 Net with different correlation matrices on P12 488 to verify the effectiveness.

Table 6: Ablation studies on different strategies of MuSiCNet in classification.  $\checkmark(\times)$  indicates the component has (not) been applied.

490									
491	Corr Matrix		AUCROC	AUPRC		Component		P	12
492	Ones		$66.7 \pm 2.2$	$25.2 \pm 0.3$	Corr Matrix	Adjustment	Contrastive	AUROC	AUPRC
493	Rand Diag		$\begin{array}{c} 84.7 \pm 0.8 \\ 84.2 \pm 0.8 \end{array}$	$\begin{array}{c}52.2 \pm 3.2 \\48.2 \pm 3.4\end{array}$	$\checkmark$	$\checkmark$	$\checkmark$	<b>86.1</b> $\pm$ 0.4	$54.1 \pm 2.2$
494	I-GAK		$\frac{85.1}{21.0} \pm 0.6$	$\frac{52.8}{46.9} \pm 3.0$	×	$\checkmark$	$\checkmark$	$85.5 \pm 0.3$	$\underline{53.0} \pm 2.1$
495	I-DT W	<u> </u>	$81.9 \pm 0.6$	40.9±3.0	√ .(	×	×	$85.2 \pm 0.6$ $85.4 \pm 0.4$	$52.6 \pm 2.5$ 53.0 ± 2.5
496	LSP-DTW		$86.1 \pm 0.4$	$54.1 \pm 2.2$	<b>√</b>	$\hat{\checkmark}$	×	$85.4 \pm 0.4$	$52.9 \pm 2.8$
497					×	×	×	$84.2 \pm 0.8$	$48.2 \pm \textbf{3.4}$

#### ABLATION ANALYSIS AND EFFICIENCY EVALUATION 4.5

Taking P12 in the classification task with a batch size of 50 as an example, we conduct the ablation 502 study to assess the necessity of two fundamental components of MuSiCNet: correlation matrix and multi-scale learning reflected in reconstruction results adjustment and contrastive learning. As shown 504 in Table 6, the complete MuSiCNet framework, incorporating all components, achieves the best 505 performance. The absence of any component leads to varying degrees of performance degradation, as 506 evidenced in layers two to five. The second layer, which retains the multi-scale learning, exhibits 507 the second-best performance, underscoring the critical role of multi-scale learning in capturing 508 varied temporal dependencies and enhancing feature extraction. Conversely, the version lacking all 509 components shows a significant performance drop of 1.9%, indicating that each component is crucial 510 to the overall effectiveness of the framework.

511 Under the same setting, our model MuSiCNet achieves a time cost of 0.240s per batch with 4.2GB of 512 memory usage. In comparison, ViTST requires 2.196s and 40.2GB, Raindrop uses 0.124s and 4.8GB, 513 MTGNN takes 0.1967s and 4.2GB, and DGM<sup>2</sup>-O needs 0.313s and 9.1GB. MuSiCNet demonstrates 514 lower time complexity than most other methods and significantly lower memory usage, particularly 515 compared to ViTST, which also performs well on classification tasks. 516

517

489

498 499 500

501

#### 5 CONCLUSION

518 519

In this study, we introduce MuSiCNet, an innovative framework designed for analyzing ISMTS 521 datasets. MuSiCNet addresses the challenges arising from data irregularities and shows superior 522 performance in both supervised and unsupervised tasks. We recognize that irregularities in ISMTS are 523 inherently relative and accordingly implement multi-scale learning, a vital element of our framework. In this multi-scale approach, the contribution of extra coarse-grained relatively regular series is 524 important, providing comprehensive temporal insights that facilitate the analysis of finer-grained 525 series. As another key component of MuSiCNet, CorrNet is engineered to aggregate temporal 526 information effectively, employing time embeddings and correlation matrix calculating from both 527 intra- and inter-series perspectives, in which we employ LSP-DTW to develop frequency correlation 528 matrices that not only reduce the burden for similarity calculation for ISMT, but also significantly 529 enhance inter-series information extraction. 530

Nevertheless, our MuSiCNet still has some limitations. Firstly, the interaction between scales could 531 potentially be further simplified. Secondly, the exploration of ISMTS for anomaly detection tasks is 532 currently insufficient. As a task-agnostic model, our MuSiCNet should be further investigated for its 533 potential in anomaly detection. 534

535

#### 536 REFERENCES 537

Rohit Agarwal, Aman Sinha, Dilip K Prasad, Marianne Clausel, Alexander Horsch, Mathieu Constant, 538 and Xavier Coubez. Modelling irregularly sampled time series without imputation. arXiv preprint arXiv:2309.08698, 2023.

540	Inci M Baytas, Cao Xiao, Xi Zhang, Fei Wang, Anil K Jain, and Jiayu Zhou. Patient subtyning via							
541	time-aware lstm networks. In Proceedings of the 23rd ACM SIGKDD international conference on							
542	knowledge discovery and data mining, pp. 65–74, 2017.							
543								
544	Donald J Berndt and James Clifford. Using dynamic time warping to find patterns in time series.							
545	In Proceedings of the 3rd international conference on knowledge discovery and data mining, pp. 250, 270, 1004							
546	557-570, 1994.							
547	Marin Biloš, Johanna Sommer, Syama Sundar Rangapuram, Tim Januschowski, and Stephan Gün-							
548	nemann. Neural flows: Efficient alternative to neural odes. Advances in neural information							
549	processing systems, 34:21325–21337, 2021.							
550	Wanlin Cai, Yuxuan Liang, Xianggen Liu, Jianshuai Feng, and Yuankai Wu, Msgnet: Learning							
551	multi-scale inter-series correlations for multivariate time series forecasting. In <i>Proceedings of the</i>							
552	AAAI Conference on Artificial Intelligence, volume 38, pp. 11141–11149, 2024.							
553	Demine D. Camine Christian A. Hammansharitk and Dada State Incompany missing data incompany							
554	with deep generative models arXiv preprint arXiv:1902.10666, 2019							
556	with deep generative models. arXiv preprint arXiv.1902.10000, 2019.							
557	Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural							
558	networks for multivariate time series with missing values. <i>Scientific reports</i> , 8(1):1–12, 2018.							
559	Ricky TO Chen Yulia Rubanova Jesse Bettencourt and David K Duvenaud Neural ordinary							
560	differential equations. <i>NeurIPS</i> , 31, 2018.							
561								
562	Xinyu Chen, Chengyuan Zhang, Xi-Le Zhao, Nicolas Saunier, and Lijun Sun. Nonstationary temporal							
563	matrix factorization for multivariate time series forecasting. <i>arXiv preprint arXiv:2203.10651</i> , 2022							
564	2022.							
565	Zipeng Chen, Qianli Ma, and Zhenxi Lin. Time-aware multi-scale rnns for time series modeling. In							
566	<i>IJCAI</i> , pp. 2285–2291, 2021.							
567	Ranak Roy Chowdhury Jiacheng Li Xiyuan Zhang Dezhi Hong Rajesh K Gunta and Jingbo							
568	Shang, Primenet: Pre-training for irregular multivariate time series. In <i>Proceedings of the AAAI</i>							
569	Conference on Artificial Intelligence, volume 37, pp. 7184–7192, 2023.							
570								
571	Marco Cuturi. Fast global alignment kernels. In <i>Proceedings of the 28th international conference on</i>							
572	<i>machine learning (ICML-11)</i> , pp. 929–930, 2011.							
573	Edward De Brouwer, Jaak Simm, Adam Arany, and Yves Moreau. Gru-ode-bayes: Continuous							
574	modeling of sporadically-observed time series. NeurIPS, 32, 2019.							
575	Wenije Du David Côté and Yan Liu Saits: Self-attention-based imputation for time series <i>Franct</i>							
576	Systems with Applications, 219:119619, 2023.							
577								
570	Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and							
580	unristoph Feichtenhoter. Multiscale vision transformers. In <i>Proceedings of the IEEE/CVF</i>							
581	international conjerence on computer vision, pp. 0824–0855, 2021.							
582	Jicong Fan. Dynamic nonlinear matrix completion for time-varying data imputation. In AAAI, March							
583	2022.							
584	Ary I. Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Poger G							
585	Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank.							
586	physiotoolkit, and physionet: components of a new research resource for complex physiologic							
587	signals. <i>circulation</i> , 101(23):e215–e220, 2000.							
588	Max Horn Michael Moor Christian Deak Destion Disak and Version Densmandt Set for sting for							
589	time series. In International Conference on Machine Learning, pp. 4353_4363 PMLP 2020							
590	and series. In microautonal conjectnice on machine Learning, pp. 4555–4505. I WER, 2020.							
591	Sheo Yon Jhin, Jaehoon Lee, Minju Jo, Seungji Kook, Jinsung Jeon, Jihyeon Hyeong, Jayoung							
592	Kim, and Noseong Park. Exit: Extrapolation and interpolation-based neural controlled differential							
593	equations for time-series classification and forecasting. In <i>Proceedings of the ACM Web Conference</i> 2022, pp. 3102–3112, 2022.							

- 594 Ming Jin, Yu Zheng, Yuan-Fang Li, Siheng Chen, Bin Yang, and Shirui Pan. Multivariate time 595 series forecasting with dynamic graph neural odes. IEEE Transactions on Knowledge and Data 596 Engineering, 2022. 597 AE Johnson, Tom J Pollard, Lu Shen, L-w H Lehman, Mengling Feng, Mohammad Ghassemi, 598 Benjamin Moody, Peter Szolovits, L Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database sci. Data, 3(1):1, 2016. 600 601 Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay 602 Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus Brubaker. Time2vec: Learning a 603 vector representation of time. arXiv preprint arXiv:1907.05321, 2019. 604 Patrick Kidger, James Morrill, James Foster, and Terry Lyons. Neural controlled differential equations 605 for irregular time series. NeurIPS, 33:6696-6707, 2020. 606 607 Zekun Li, Shiyang Li, and Xifeng Yan. Time series as images: Vision transformer for irregularly 608 sampled time series. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https://openreview.net/forum?id=ZmeAoWQqe0. 609 610 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dust-611 dar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and 612 forecasting. In ICLR, 2021. 613 Nicholas R Lomb. Least-squares frequency analysis of unequally spaced data. Astrophysics and 614 space science, 39:447-462, 1976. 615 616 Stephane G Mallat. A theory for multiresolution signal decomposition: the wavelet representation. 617 *IEEE transactions on pattern analysis and machine intelligence*, 11(7):674–693, 1989. 618 Matthew J Menne, CN Williams Jr, and Russell S Vose. United states historical climatology network 619 daily temperature, precipitation, and snow data. Carbon Dioxide Information Analysis Center, Oak 620 Ridge National Laboratory, Oak Ridge, Tennessee, 2015. 621 622 Piotr Nawrot, Szymon Tworkowski, Michał Tyrolski, Łukasz Kaiser, Yuhuai Wu, Christian Szegedy, 623 and Henryk Michalewski. Hierarchical transformers are more efficient language models. arXiv 624 preprint arXiv:2110.13711, 2021. 625 Daniel Neil, Michael Pfeiffer, and Shih-Chii Liu. Phased lstm: Accelerating recurrent network 626 training for long or event-based sequences. *NeurIPS*, 29, 2016. 627 628 Attila Reiss and Didier Stricker. Introducing a new benchmarked dataset for activity monitoring. In 629 2012 16th international symposium on wearable computers, pp. 108–109. IEEE, 2012. 630 Matthew A Reyna, Christopher S Josef, Russell Jeter, Supreeth P Shashikumar, M Brandon Westover, 631 Shamim Nemati, Gari D Clifford, and Ashish Sharma. Early prediction of sepsis from clinical data: 632 the physionet/computing in cardiology challenge 2019. Critical care medicine, 48(2):210-217, 633 2020. 634 Yulia Rubanova, Ricky TQ Chen, and David K Duvenaud. Latent ordinary differential equations for 635 irregularly-sampled time series. NeurIPS, 32, 2019. 636 637 Jeffrey D Scargle. Studies in astronomical time series analysis. ii-statistical aspects of spectral 638 analysis of unevenly spaced data. Astrophysical Journal, Part 1, vol. 263, Dec. 15, 1982, p. 639 835-853., 263:835-853, 1982. 640 Mona Schirmer, Mazin Eltayeb, Stefan Lessmann, and Maja Rudolph. Modeling irregular time 641 series with continuous recurrent units. In International Conference on Machine Learning, pp. 642 19388-19405. PMLR, 2022. 643 644 Mohammad Amin Shabani, Amir H Abdi, Lili Meng, and Tristan Sylvain. Scaleformer: Iterative 645 multi-scale refining transformers for time series forecasting. In The Eleventh ICLR, 2022. 646
- 647 Satya Narayan Shukla and Benjamin Marlin. Interpolation-prediction networks for irregularly sampled time series. In *ICLR*, 2018.

648 Satya Narayan Shukla and Benjamin Marlin. Multi-time attention networks for irregularly sampled 649 time series. In ICLR, 2021. URL https://openreview.net/forum?id=4c0J6lwQ4\_. 650 Satya Narayan Shukla and Benjamin M Marlin. A survey on principles, models and methods for 651 learning from irregularly sampled time series. arXiv preprint arXiv:2012.00168, 2020. 652 653 Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-hospital 654 mortality of icu patients: The physionet/computing in cardiology challenge 2012. In 2012 655 Computing in Cardiology, pp. 245–248. IEEE, 2012. 656 657 Bhanu Pratap Singh, Iman Deznabi, Bharath Narasimhan, Bryon Kucharski, Rheeya Uppaal, Akhila Josyula, and Madalina Fiterau. Multi-resolution networks for flexible irregular time series modeling 658 (multi-fit). arXiv preprint arXiv:1905.00125, 2019. 659 660 Chenxi Sun, Shenda Hong, Moxian Song, Yen-Hsiu Chou, Yongyue Sun, Derun Cai, and Hongyan 661 Li. Te-esn: Time encoding echo state network for prediction based on irregularly sampled time 662 series data. In Zhi-Hua Zhou (ed.), IJCAI, pp. 3010–3016. International Joint Conferences on 663 Artificial Intelligence Organization, 8 2021. 665 Chenxi Sun, Hongyan Li, Moxian Song, Derun Cai, Baofeng Zhang, and Shenda Hong. Time pattern reconstruction for classification of irregularly sampled time series. Pattern Recognition, 147: 666 110075, 2024. 667 668 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csdi: Conditional score-based 669 diffusion models for probabilistic time series imputation. NeurIPS, 34, 2021. 670 671 Jacob T VanderPlas. Understanding the lomb-scargle periodogram. The Astrophysical Journal Supplement Series, 236(1):16, 2018. 672 673 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 674 Kaiser, and Illia Polosukhin. Attention is all you need. NeurIPS, 30, 2017. 675 676 Cédric Villani et al. Optimal transport: old and new, volume 338. Springer, 2009. 677 Jun Wang, Wenjie Du, Wei Cao, Keli Zhang, Wenjia Wang, Yuxuan Liang, and Qingsong Wen. Deep 678 learning for multivariate time series imputation: A survey. arXiv preprint arXiv:2402.04059, 2024. 679 680 Yinjun Wu, Jingchao Ni, Wei Cheng, Bo Zong, Dongjin Song, Zhengzhang Chen, Yanchi Liu, 681 Xuchao Zhang, Haifeng Chen, and Susan B Davidson. Dynamic gaussian mixture based deep 682 generative model for robust forecasting on sparse multivariate time series. In AAAI, volume 35, pp. 683 651-659, 2021. 684 Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Connecting 685 the dots: Multivariate time series forecasting with graph neural networks. In Proceedings of the 686 26th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 753–763, 687 2020. 688 689 Vijaya Krishna Yalavarthi, Kiran Madhusudhanan, Randolf Scholz, Nourhan Ahmed, Johannes 690 Burchert, Shayan Jawed, Stefan Born, and Lars Schmidt-Thieme. Grafiti: Graphs for forecasting 691 irregularly sampled time series. In Proceedings of the AAAI Conference on Artificial Intelligence 692 (AAAI), pp. 16255–16263, 2024. 693 Jinsung Yoon, James Jordon, and Mihaela Schaar. Gain: Missing data imputation using generative 694 adversarial nets. In International conference on machine learning, pp. 5689–5698. PMLR, 2018. 696 Zhihao Yu, Xu Chu, Liantao Ma, Yasha Wang, and Wenwu Zhu. Imputation with inter-series 697 information from prototypes for irregular sampled time series. arXiv preprint arXiv:2401.07249, 2024. 699 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series 700 forecasting? In Proceedings of the AAAI conference on artificial intelligence, volume 37, pp.

11121-11128, 2023.

- 702 Jiawen Zhang, Shun Zheng, Wei Cao, Jiang Bian, and Jia Li. Warpformer: A multi-scale modeling 703 approach for irregular clinical time series. In Proceedings of the 29th ACM SIGKDD Conference 704 on Knowledge Discovery and Data Mining, pp. 3273–3285, 2023. 705 706 Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. 707 Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In Proceedings of the IEEE/CVF international conference on computer vision, pp. 2998–3008, 2021a. 708 709 Xiang Zhang, Marko Zeman, Theodoros Tsiligkaridis, and Marinka Zitnik. Graph-guided network 710 for irregularly sampled multivariate time series. In ICLR, 2021b. 711 712 Zhao-Yu Zhang, Shao-Qun Zhang, Yuan Jiang, and Zhi-Hua Zhou. Life: Learning individual features 713 for multivariate time series prediction with missing values. In 2021 IEEE International Conference 714 on Data Mining (ICDM), pp. 1511–1516. IEEE, 2021c. 715 716 Yucheng Zhao, Chong Luo, Zheng-Jun Zha, and Wenjun Zeng. Multi-scale group transformer for 717 long sequence modeling in speech separation. In IJCAI, pp. 3251–3257, 2021. 718 719 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 720 Informer: Beyond efficient transformer for long sequence time-series forecasting. In Proceedings 721 of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021. 722 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency 723 enhanced decomposed transformer for long-term series forecasting. In International conference on 724 machine learning, pp. 27268-27286. PMLR, 2022. 725 726 727
- 728 729

730

731 732

733

734

735

736 737

738 739

740

741

А APPENDIX

#### R PSEUDO CODE FOR MUSICNET

The Pseudo Code is provided using classification as an example. The interpolation task can be obtained by removing the projection head  $f_{cls}$  and the classification loss term  $\mathcal{L}_{cls}$  from the total loss in line #17. While in the case of forecasting tasks, the projection head will be replaced with  $f_{\text{fore}}$  and task loss will be changed to  $\mathcal{L}_{\text{fore}}$  as in Eq.11.

#### С TIME EMBEDDING IN CORRNET

Time Embedding method embeds continuous time points of ISMTS into a vector space Kazemi et al. (2019); Shukla & Marlin (2021). It leverages H embedding functions  $\phi_h(t)$  simultaneously and each outputting a representation of size  $d_r$ . Dimension *i* of embedding *h* is defined as follows:

$$\phi_h(t)[i] = \begin{cases} \omega_{0h} \cdot t + \alpha_{0h}, & \text{if } i = 0\\ \sin\left(\omega_{ih} \cdot t + \alpha_{ih}\right), & \text{if } 0 < i < d_r \end{cases}$$
(7)

746 where the  $\omega_{ih}$ 's and  $\alpha_{ih}$ 's are learnable parameters that represent the frequency and phase of the sine function. This time embedding method can capture both non-periodic and periodic patterns with 747 linear and periodic terms, respectively. c 748

749 750 751

752

753 754 755

#### D **ISMTS ANALYSIS TASKS**

The overall loss is defined as Eq.8, incorporating an optional task-specific loss component.

$$\mathcal{L} = \sum_{l=1}^{L} \ell_{\text{recon}}^{(l)} + \lambda_1 \ell_{\text{adj}}^{(l)} + \lambda_2 \ell_{\text{cons}}^{(l)}$$
(8)

756 Algorithm 1 MuSiCNet Algorithm for Classification 757 **Input:** Training set  $\mathcal{X}$ , the number of scale layers L, random masking ratio r, max reference point 758 number  $|\boldsymbol{\tau}^{(L)}|$ , hyper-parameters  $\lambda_1, \lambda_2, \lambda_3$ . 759 **Parameters:** Encoder model  $f_{\text{CorrE}}$ , decoder model  $f_{\text{CorrD}}$ , GRU model  $f_{\text{GRU}}$ , projection head  $f_{\text{cls}}$ 760 **Output:** Encoder model  $f_{\text{CorrE}}$ , GRU model  $f_{\text{GRU}}$ , projection head  $f_{\text{cls}}$ 761 1:  $C_T \leftarrow \text{Eq.2}$  with  $\mathcal{X}$ 762 2: for X in  $\mathcal{X}$  do  $\left\{ \boldsymbol{X}^{(1)}, \cdots, \boldsymbol{X}^{(L)} \right\} \leftarrow \operatorname{Mask}_{r}\left(\operatorname{AvgPooling}_{L}\left(\boldsymbol{X}\right)\right)$ 763 3: 764 4:  $\ell_{\text{recon}} \leftarrow 0$ 5: 765 for  $l \leftarrow 1$  to L do  $\boldsymbol{r}^{(l)} \leftarrow f_{\text{CorrE}}\left(X^{(L)}, C_T, |\boldsymbol{\tau}^{(L)}|/2^{(L-l)}\right)$ 766 6:  $\boldsymbol{h}^{(l)} \leftarrow f_{\text{GRU}}\left(\boldsymbol{r}^{(l)}\right)$ 767 7: 768  $\hat{\boldsymbol{X}}_{\text{recon}}^{(l)} \leftarrow f_{\text{CorrD}}\left(\boldsymbol{r}^{(l)}, |\boldsymbol{X}^{(l)}|\right)$ 8: 769  $\ell_{\text{recon}} \leftarrow \ell_{\text{recon}} + \text{Eq.3} \text{ with } \boldsymbol{X}^{(l)} \text{ and } \hat{\boldsymbol{X}}^{(l)}$ 9: 770 end for 10: 771 11:  $\ell_{adj}, \ell_{cons} \leftarrow 0, 0$ 772 for  $l \leftarrow 2$  to L do 12:  $\ell_{\text{adj}} \leftarrow \ell_{\text{adj}} + \text{Eq.5} \text{ with } \hat{m{X}}^{(l-1)} \text{ and } \hat{m{X}}^{(l)}$ 773 13: 774  $\ell_{cons} \leftarrow \ell_{cons} + \text{Eq.6} \text{ with } \boldsymbol{h}^{(l-1)} \text{ and } \boldsymbol{h}^{(l)}$ 14: 775 end for 15:  $\mathcal{L}_{cls} \leftarrow Eq.9$  with  $h^{(L)}$ 776 16:  $\mathcal{L}_{\text{overall}} \leftarrow \frac{1}{L} \ell_{\text{recon}} + \frac{\lambda_1}{L-1} \ell_{\text{adj}} + \frac{\lambda_2}{L-1} \ell_{\text{cons}} + \lambda_3 \mathcal{L}_{\text{cls}}$ 777 17: 778 18: Update overall network parameters 779 19: end for

780 781

791 792

793

794

796 797 798

799 800

801

802

803 804 805

Supervised Learning. We augment the encoder-decoder CorrNet by integrating a supervised
 learning component that utilizes the latent representations for feature extraction. In this work, we
 specifically concentrate on classification tasks as a representative example of supervised learning.
 The loss function is

$$\mathcal{L}_{\text{cls}} = \frac{1}{C} \sum_{c=1}^{C} \frac{1}{n^c} \sum_{i=1}^{n^c} \ell_{CE} \left( \text{CLS} \left( \boldsymbol{h}_i^{(L)} \right), y_i \right)$$
(9)

where C denotes the number of classes,  $n^c$  denotes the number of samples in c-th class, CLS (·) denotes the projection head for classification, and  $\ell_{CE}$  (·) denotes the cross-entropy loss.

**Unsupervised Learning.** For our unsupervised learning example, we choose interpolation and forecasting. The loss function for interpolation is defined as

$$\mathcal{L}_{int} = \sum_{n=1}^{N} \left\| \boldsymbol{M}^{(L)} \odot \left( (\hat{\boldsymbol{X}}_{\text{reco}}^{(L)})_n - \boldsymbol{X}_n^{(L)} \right) \right\|_2^2$$
(10)

This equation essentially represents the reconstruction outcome at the finest scale as  $\ell_{adj}^{(L)}$  in Eq.4 making the interpolation task fit seamlessly into our model with minimal modifications. Therefore, it is unnecessary to incorporate an additional loss function into our overall loss function Eq.8.

While the loss function for forecasting is defined as

$$\mathcal{L}_{\text{fore}} = \sum_{n=1}^{N} \left\| (\boldsymbol{M}_{\text{fore}})_n \odot \left( (\hat{\boldsymbol{X}}_{\text{fore}}^{(L)})_n - (\boldsymbol{X}_{\text{fore}})_n \right) \right\|_2^2$$
(11)

806 807 808

As observations might be missing also in the groundtruth data, to measure forecasting accuracy we average an element-wise loss function  $\mathcal{L}_{\text{fore}}$  over only valid values using  $(M_{\text{fore}})_n$ .

Tasks	Datasets	#Samples	#Variables	#Avg. obs.	#Classes	Imbalanced	Missing ratio
Classification	P19	38,803	34	401	2	True	94.9%
	P12	11,988	36	233	2	True	88.4%
	PAM	5,333	17	4,048	8	False	60.0%
Interpolation	PhysioNet	4,000	37	2,880	-	-	78.0%
Forecasting	USHCN	1,100	5	263	-	-	77.9%
	MIMIC-III	21,000	96	274	-	-	94.2%
	Physionet12	5,333	37	130	-	-	85.7%

Table 7: Statistics of the ISMTS datasets used in our experiments. "#Avg. obs." denotes the average
 number of observations for each sample.

## E FURTHER DETAILS ON DATASETS

We adopt the data processing approach used in RAINDROP Zhang et al. (2021b) for the classification task, mTANs Shukla & Marlin (2021) for the interpolation task, and GraFITi Yalavarthi et al. (2024) for the forecasting task. The aforementioned processing methods serve as the usual setup, which our method also follows for fair comparison. *However, it's important to note that we do not incorporate static attribute vectors* (such as age, gender, time from hospital to ICU admission, ICU type, and length of stay in ICU) in our processing. This decision is based on the fact that our model, MuSiCNet, is not specifically designed for clinical datasets. Instead, it is designed as a versatile, general model capable of handling various types of datasets, which may not always include such static vectors. The detailed information of baselines is in Table 7.

833 E.1 DATASETS FOR CLASSIFICATION

P19: PhysioNet Sepsis Early Prediction Challenge 2019. P19 dataset Reyna et al. (2020) comprises data from 38, 803 patients, each monitored by 34 irregularly sampled sensors, including 8 vital signs and 26 laboratory values. The original dataset contained 40, 336 patients, but we excluded those with excessively short or long time series, resulting in a range of 1 to 60 observations per patient as in RAINDROP. Each patient has a binary label representing the occurrence of sepsis within the next 6 hours. The dataset has a high imbalance with approximately ~ 4% positive samples.

841

847

822 823

824

825

826

827

828

829

830

831

832

834

P12: PhysioNet Mortality Prediction Challenge 2012. P12 Goldberger et al. (2000) includes data from 11, 988 patients after removing inappropriate 12 samples as explained in Horn et al. (2020). This dataset features multivariate time series from 36 sensors collected during the first 48 hours of ICU stay. Each patient has a binary label indicating the length of stay in the ICU, in which a negative label for stays under 3 days and a positive label for longer stays. P12 is imbalanced with ~ 93% positive samples.

PAM: PAMAP2 Physical Activity Monitoring. PAM Reiss & Stricker (2012) records the daily 848 activities of 9 subjects using 3 inertial measurement units. RAINDROP has adapted it for irregularly 849 sampled time series classification by excluding the ninth subject for short sensor data length. The 850 continuous signals were segmented into samples with the window size 600 and 50% overlapping 851 rate. Originally with 18 activities, we retain 8 with over 500 samples each, while others are dropped. 852 After modification, PAM includes 5, 333 sensory signal segments, each with 600 observations from 853 17 sensors at 100 Hz. To simulate irregularity, 60% of observations are randomly removed by 854 RAINDROP, uniformly across all experimental setups for fair comparison. The 8 classes of PAM 855 represent different daily activities, with no static attributes and roughly balanced distribution.

856 857

858

E.2 DATASET FOR INTERPOLATION

Physionet: PhysioNet Challenge 2012 dataset Physionet Reiss & Stricker (2012) comprises 37
variables from ICU patient records, with each record containing data from the first 48 hours after
admission to ICU. Aligning with the methodology of Neural ODE Rubanova et al. (2019), we round
observation times to the nearest minute, resulting in up to 2, 880 potential measurement times for
each time series. The dataset encompasses 4, 000 labeled instances and an equal number of unlabeled
instances. For our study, we utilize all 8, 000 instances in interpolation experiments. Our primary

objective is to predict in-hospital mortality, with 13.8% of the instances belonging to the positive class.

866 867 868

## E.3 DATASET FOR FORECASTING

USHCN: U.S. Historical Climatology Network. USHCN Menne et al. (2015) data are used to quantify national and regional-scale temperature changes in the contiguous United States. It contains measurements of 5 variables from 1280 weather stations. Following the preprocessing proposed by De Brouwer et al. (2019), the majority of the over 150 years of observations are excluded, and only data from the years 1996 to 2000 are used in the experiments. Furthermore, to create a sparse dataset, only a randomly sampled 5% of the measurements are retained.

875

Physionet12. This dataset consists of medical records from 12,000 ICU patients. During the first
48 hours of admission, measurements of 37 vital signs were recorded. Following the forecasting
approach used in recent work, such as Yalavarthi et al. (2024); Biloš et al. (2021); De Brouwer
et al. (2019), we pre-process the dataset to create hourly observations, resulting in a maximum of 48
observations per series.

- MIMIC-III: Medical Information Mart for Intensive Care. MIMIC-III Johnson et al. (2016)
  is a widely utilized medical dataset offering valuable insights into ICU patient care. To capture
  a diverse range of patient characteristics and medical conditions, 96 variables are meticulously
  observed and documented. For consistency, we followed the preprocessing steps outlined in previous
  studies Yalavarthi et al. (2024); Schirmer et al. (2022); Biloš et al. (2021); De Brouwer et al. (2019).
  Specifically, we rounded the recorded observations to 30-minute intervals and used only the data from
  the first 48 hours post-admission. Patients who spent less than 48 hours in the ICU were excluded
  from the analysis.
  - 889 890

## F EXPERIMENTAL DETAILS

891 892 893

## F.1 MUSICNET PARAMETERS

We present the training hyperparameters and model parameters here. The maximum epoch is set to 300, and AdamW optimizer is selected as our optimizer without weight decay. By default, the learning rate is set to 1e-3, and the learning rate schedule is cosine decay for each epoch. Batch size for all datasets is set to 50, the dimension of the encoder output is set to 256, and the dimension of the hidden representations in GRU is typically set to 50. The random masking ratio r for each scale is set to 0.1.

Due to inconsistent series lengths, we set the maximum reference point number, K, to 128 for long series, such as P12, PAM, PhysioNet and USHCN, to 96 for Physionet12, and to 48 for short series, such as PAM and MIMIC-III.

Initially, the window size is set to 1/4 of the time series length and then halved iteratively until the majority of the windows contain at least one observation.

According to the observed timestamps on each dataset, the number of scale layers L is set to 6, 5, 7, 6, 8, 4, and 5 for P12, P19, PAM, Physionet, USHCN, MIMIC-III, and Physionet12, respectively. For example, in classification, for P12, the scales are 4, 8, 16, 32, 64 and raw length. For P19, the scales are 4, 8, 16, 32 and raw length. And for PAM, the scales are 4, 8, 16, 32, 64, 128 and raw length. In all mainstream tasks involved, the hyperparametes  $\lambda_1, \lambda_2, \lambda_3$  are selected in [1e-3, 1e-2, ..., 1e2]. All the models were experimented using the PyTorch library on a GeForce RTX-2080 Ti GPU.

912

# 913 F.2 BASELINE PARAMETERS

The implementation of baseline models adheres closely to the methodologies outlined in their respective papers, including SeFT Horn et al. (2020), GRU-D Che et al. (2018), mTAND Shukla & Marlin (2021) and ViTST Li et al. (2023). We follow the settings of the attention embedding module baseline in mTAND and implement the Multi-Correlation Attention module in our work.





Figure 4: AUCROC performance with varying Figure 5: AUCROC performance with varying combinations of hyper-parameter of the adjust-hyper-parameter of the downstream task  $\lambda_3$  in the ment term  $\lambda_1$  and hyper-parameter of the con-logarithmic form on P12 trastive learning term  $\lambda_2$  in the logarithmic form on P12

#### F.3 PARAMETER ANALYSIS

To analyze the hyper-parameters sensitivity, we conducted the experiments for  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  with grid search. Due to the closer relationship between the hyper-parameters of the adjustment term and the contrastive learning term, i.e.,  $\lambda_1$  and  $\lambda_2$ , we jointly analyzed  $\lambda_1$  and  $\lambda_2$  while separately analyzing the hyper-parameter of the downstream task  $\lambda_3$ , as illustrated in Fig.4 and Fig.5. 

From Fig.4, we can find that the adjustment term plays a greater role compared to the contrastive learning term. This phenomenon matches our motivation, where the coarse-to-fine strategy can effectively alleviate the difficulty of representation learning on ISMTS caused by inconsistent sampling rates. In addition, when  $\lg \lambda_1$  and  $\lg \lambda_2$  take values around 2 and -2, respectively, our MuSiCNet can perform well. 

From Fig.5, we can find that our MuSiCNet becomes effective with large  $\lambda_3$ . This indicates that more effective representations will be captured when utilizing downstream tasks, matching the general insight. We also noticed that it becomes less sensitive when  $\lg \lambda_3 \ge -1$ . Its suitable range may be located at [1e1, 1e2].