# FEDNET: FREQUENCY ENHANCED DECOMPOSED NETWORK FOR OUT-OF-DISTRIBUTION TIME SERIES CLASSIFICATION

Anonymous authors

Paper under double-blind review

## Abstract

Time series classification is a crucial task with widespread applications in various fields such as medicine and energy. Due to the non-stationary property of time series, its data distribution will change over time, which makes it challenging for models to generalize to the *out-of-distribution* (OOD) environment. However, limitations persist in the current research on OOD time series classification, particularly the absence of a unified consideration addressing both domain distribution shift and temporal distribution shift. To this end, we view the time series distribution shift from the frequency perspective and propose a novel method called *Frequency Enhanced Decomposed Network* (FEDNet) for OOD time series classification. FEDNet utilizes frequency domain information to guide the decomposition of time series and further eliminates domain shift and temporal shift, it then obtains domain-invariant features for adapting to OOD data. Finally, we provide theoretical insights of FEDNet to validate its superiority for OOD time series classification. Comprehensive results on synthetic and real-world datasets demonstrate that FEDNet achieves state-of-the-art performance in OOD time series classification tasks, surpassing previous methods by up to 7%. Our code is available at https://anonymous.4open.science/r/FEDNet-743E.

028 029

031

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

## 1 INTRODUCTION

Time series classification is a pivotal task with applications in bio-signals processing (Salehi et al., 2021), medical diagnostics (Supratak et al., 2017), and human activity recognition (Tang et al., 2020). Recent studies have introduced methods like TCN (Bai et al., 2018) and shallow RNN (Dennis et al., 2019) for this purpose, largely adhering to the *independently identically distributed* (i.i.d) assumption. However, this assumption no longer holds in reality as the testing data do not always follow the same distribution as the training data, i.e. *out-of-distribution* (OOD), thus the performance on the testing data is severely degraded.

In practical situations, it proves to be highly challenging to acquire the distribution of the testing set 040 (Wang et al., 2022b). As shown in Figure 1, we can only get training samples from a finite number 041 of domains, these domains usually represent different types of populations, while the testing data is 042 invisible and inaccessible in reality. Therefore, we cannot utilize the testing data for domain adap-043 tation (Patel et al., 2015) and the limited samples are not enough to build a powerful pre-training 044 model for transfer learning (Pan & Yang, 2009). How to use limited domain datasets to improve the model's generalization on datasets in unseen domains becomes a realistic problem. Moreover, traditional domain generalization methods following the invariant risk minimization (IRM) paradigm 046 are suboptimal to OOD time series classification since the marginal probability distribution of time 047 series data would change when facing non-stationary situations, which may cause their distribution 048 to deviate from the corresponding domains. 049

Unfortunately, research on OOD generalization for time series classification is limited. Existing
methods can be divided into two main types, domain relabeling and disentangled representation
learning methods. On the one hand, the domain relabeling methods (Lu et al., 2022; Du et al., 2021)
try to divide the time series into groups of segments with large distance in data distribution and
relabel its domain for training. On the other hand, the disentangled representation learning methods



Figure 1: OOD time series classification scenario. The cross-people division (by age or other rules) causes domain distribution shift, while different periods caused temporal distribution shift.

(Qian et al., 2021) expect to separate the data into domain-invariant and domain-specific parts. The domain-invariant part represents features that remain consistent or unchanged across domains, while the domain-specific part represents features that exhibit variations or specificity in different domains. Nonetheless, domain relabeling methods directly utilize the whole segments to reset domain labels but ignore the noise and redundant information between segments, and disentangled representation learning methods ignore temporal distribution shift within the domain by directly utilizing the domain labels to aid training. In addition, these methods analyze distribution shift directly from the time domain without considering the global view of the frequency information in time series.

To fill the gap of frequency information and deal with temporal distribution shift in OOD time series 077 classification, we propose a method called *Frequency Enhanced Decomposed Network* (FEDNet). Inspired by Wold's Theorem (Jenkins et al., 1955) and Koopa's (Liu et al., 2023) successful prac-079 tice in non-stationary time series forecast, we realized that time series can be decomposed into time-deterministic component and time-stochastic component by frequency information. The time-081 deterministic part is less affected by the temporal changes, while the time-stochastic part mainly 082 depends on temporal changes. It encouraged us to introduce frequency domain information and 083 time series decomposition ideas into OOD time series classification to analyze the data in differ-084 ent domains. Moreover, we use both fast Fourier transform (FFT) and Discrete Wavelet Trans-085 form (DWT) to extract a certain percentage of high average amplitude frequency components as time-deterministic features by scanning all training data in different domains. Specifically, the timedeterministic features are less affected by time and can be decoupled to extract domain-invariant 087 880 features, and the remaining frequency components are used to model the temporal stochastic components. Finally, we theoretically compare FEDNet with previous OOD generalization methods. 089

- <sup>090</sup> The contributions of our work can be summarized as follows:
  - To the best of our knowledge, it is the first work to investigate OOD time series classification from the frequency perspective. Additionally, we postulate and formulate the concept of frequency distribution shift for modeling temporal distribution.
  - We propose a novel method called FEDNet and provide its theoretical insights. FEDNet decomposes time series into time-deterministic and time-stochastic parts separating the temporal distribution by frequency information. We propose that domain feature contrastive learning simplifies the constraints on domain-invariant learning and accelerates convergence.
  - Experiments on several datasets demonstrate the state-of-the-art performance of FEDNet and the effectiveness of frequency information. Moreover, we found the stability of the frequency component to time-distribution shifts and verified that the essence of frequency enhancement originate from the orthogonality of frequencies.
- 101 102 103

092

094

095

096

098

099

100

065

066

067

- 2 RELATED WORK
- 104 105

Domain Generalization. Domain generalization (Wang et al., 2022a) is a difficult challenge where
 the goal is to obtain robust models from multiple domains that eventually generalize to unseen
 domains and are inaccessible to the training process. Existing domain generalization schemes are

108 divided into three main levels: data augmentation, feature representation, and learning strategies. In 109 data augmentation, Mixup(Zhang, 2017) data from different domains is often used in data processing 110 and some methods try to augment data from a frequency domain perspective(Demirel & Holz, 2024; 111 Xu et al., 2021) verified the robustness of frequency information. In representation learning, a 112 common strategy is to extract domain-invariant features across multiple source domains with IRM paradigm(Arjovsky et al., 2019; Krueger et al., 2021), domain adversarial learning (Ganin et al., 113 2016) and disentanglement-based method(Ilse et al., 2020) serve the same purpose. For learning 114 strategies, one of the most famous approaches is distributional robust optimization, such methods 115 (Sagawa\* et al., 2020; Kuhn et al., 2019) dynamically penalize the weighted average loss of all 116 domain distribution sets, thus minimizing the generalization expectation of the unseen distributions 117 theoretically. Some methods attempt to decompose features to capture more stable domain-invariant 118 information. StableNet(Zhang et al., 2021), inspired by causal mechanisms, introduces a novel 119 nonlinear feature decomposition correlation technique for capturing domain invariant information. 120 Unfortunately, due to the widespread temporal distribution shift in time series, these above methods 121 cannot be fully applicable to domain generalization for time series classification.

122 **OOD Time Series Classification.** The research on OOD time series classification is very limited. 123 AdaRNN (Du et al., 2021) proposes a temporal distribution characterization and matching module 124 that divides the time series segments into finite large-distance groups in data distribution to extract 125 invariant features. Diversify (Lu et al., 2022) identifies the latent distribution domains for the "worst-126 case scenarios" through adversarial training, and then reduces the gaps between time series segments 127 in latent domains. GILE (Qian et al., 2021) is a disentanglement method designed to extract domain-128 invariant and domain-specific representations through variational inference (Kingma & Welling, 129 2013). However, these methods do not take into account both domain and temporal shifts, resulting 130 in suboptimal performance.

131 Frequency-Based Time Series Representation. Frequency domain information has been widely 132 used in recent years. In time series, many methods get better results by introducing frequency 133 domain information. For example, FEDformer (Zhou et al., 2022) improves the computational ef-134 ficiency and performance of long-short forecasting with the help of frequency attention, FreTS (Yi 135 et al., 2023) is MLP-based architecture with frequency spectrum, and TF-C (Zhang et al., 2022) pro-136 poses time-frequency consistency to do time series pre-training and transfer learning. These works 137 show that frequency domain information can improve the generalization performance of the model.

138 139

140

#### 3 PRELIMINARY

141 **Definition 1.** Time Series Data. A multi-domain time series dataset can be defined as  $\mathcal{E}$  = 142  $\{\mathcal{X}, \mathcal{Y}, \mathcal{D}\}\$ , where  $\mathcal{X}, \mathcal{Y}, \mathcal{D}$  represent time series segments pre-processed raw data samples by sliding window, category labels and domain labels respectively. Specifically,  $\mathcal{X} = \{x_i\}_{i=1}^N$  indicates N time series segments in total, where  $x_i \in \mathcal{X} \subset \mathbb{R}^{L \times C}$  is the C-channel instance with L timestamps and 143 144  $y_i \in \mathcal{Y} = \{c_1, \cdots, c_{N_y}\}$  is its label, where  $N_y$  is the number of the category labels. In addition, the 145 dataset is usually divided into multiple domains environments  $\mathcal{D} = \{d_1, d_2, \cdots, d_{N_d}\}$ , where  $N_d$  is 146 the number of the domain labels. Each domain consists of a set of samples  $D_k = \{(x_i, y_i)\}_{i=1}^{N_k}, N_k$ 147 denotes the number of samples in the k-th domain  $D_k$ . 148

149 **Definition 2.** Domain Distribution Shift. Given a multi-domain time series dataset  $\mathcal{E}$  = 150  $\{\mathcal{X}, \mathcal{Y}, \mathcal{D}\}\$ , these domains usually belong to different data distributions. We define the joint dis-151 tribution of data and label as  $\mathbb{P}(x, y)$ . Domain distribution shift can be described as 152

$$\forall D_i \neq D_j, \mathbb{P}^{D_i}(x) \neq \mathbb{P}^{D_j}(x), \mathbb{P}^{D_i}(y|x) = \mathbb{P}^{D_j}(y|x) \to \mathbb{P}^{D_i}(x,y) \neq \mathbb{P}^{D_j}(x,y).$$
(1)

154 **Definition 3. Temporal Distribution Shift.** Given a group of time series data in k-th domain 155  $D_k = \{(x_i, y_i)\}_{i=1}^{N_k}$ , suppose these segments come from  $\{D_k^t\}_{t=1}^T$  peroids, T is unpredictable. Temporal distribution shift exists between different periods caused by non-stationary or equipment factors.

153

156 157

$$\exists i, j \in [1, T], \mathbb{P}^{D_k^i}(x) \neq \mathbb{P}^{D_k^j}(x) \neq \mathbb{P}^{D_k}(x) \xrightarrow{\mathbb{P}^{D_k}(y|x)} \mathbb{P}^{D_k^i}(x, y) \neq \mathbb{P}^{D_k^j}(x, y) \neq \mathbb{P}^{D_k}(x, y).$$
(2)

**Definition 4. Frequency Distribution Shift.** The time series can be transformed to the frequency 161 domain by the Fourier transform, and we can use the frequency domain statistics to represent the



Figure 2: Architecture of the proposed FEDNet.

distribution in the same way. Here, we propose the concept of frequency distribution shift to describe the marginal probability distribution shift of segment in frequency domain.

$$\exists i, j \in [1,T], \mathbb{P}_F^{D_k^i}(x) \neq \mathbb{P}_F^{D_k^j}(x) \neq \mathbb{P}_F^{D_k}(x) \xrightarrow{\mathbb{P}_F^{D_k}(y|x)} \mathbb{P}_F^{D_k^i}(x,y) \neq \mathbb{P}_F^{D_k^j}(x,y) \neq \mathbb{P}_F^{D_k}(x,y).$$
(3)

## **PROBLEM FORMULATION**

Given a multi-domain time series dataset  $\mathcal{E} = \{\mathcal{X}, \mathcal{Y}, \mathcal{D}\}$ , we follow the cross-domain rule to divide the dataset into training dataset  $\mathcal{E}_{tr} = \{\mathcal{X}^{tr}, \mathcal{Y}^{tr}, \mathcal{D}^{tr}\}$  and unseen testing dataset  $\mathcal{E}_{te} = \{\mathcal{X}^{te}, \mathcal{Y}^{te}, \mathcal{D}^{te}\}$ . Specifically, in our paper, three conditions are imposed on the dataset: (1) OOD environments:  $\mathbb{P}^{tr}(x, y) \neq \mathbb{P}^{te}(x, y)$ . (2) Domain distribution shift:  $\mathbb{P}^{D_i}(x, y) \neq \mathbb{P}^{D_j}(x, y), \forall D_i \neq D_j$ . (3) Temporal distribution shift:  $\mathbb{P}^{D_k^i}(x, y) \neq \mathbb{P}^{D_k}(x, y), \exists i \neq j \in [1, T]$ . Our goal is to learn an optimal model  $f_{\theta}^* : \mathcal{X} \to \mathcal{Y}$  under the above conditions, the model aims to minimize both domain distribution shift and temporal distribution shift to generalize well to OOD testing data:

$$f_{\theta}^* = \arg\min_{f_{\theta}} \mathbb{E}_{(x,y)\sim\mathcal{E}_{te}}[\ell(f_{\theta}(x), y)], \tag{4}$$

where  $\mathbb{E}$  denotes expectation and  $\ell(\cdot, \cdot)$  denotes loss function.

## 5 Method

In this section, we introduce the detailed pipeline of FEDNet, which is shown in Figure 2.FEDNet decomposes time series into a time-deterministic block and a time-stochastic block to extract in-variant features for OOD time series classification. The time-deterministic block is less affected by time, its feature is mainly affected by the domain distribution shift, which can be prompted with domain labels. The time-stochastic block is more susceptible to the temporal distribution shift. We convert data into finite patches to simulate the moving average weighting process of the stochas-tic component. FEDNet decomposes the time series components theoretically to model the effects of domain distribution shift and temporal distribution shift for OOD time-series classification. Ul-timately, we theoretically analyze previous IRM-based methods and FEDNet to demonstrate that FEDNet is better suited for OOD time series classification.

# 216 5.1 FREQUENCY FILTER217

220 221

222

224

225

226

227

228

233

234

235

242 243

244

245 246

**Lemma 1 (Wold's Theorem).** Given weak-sense stationarity time series  $x_t$  can be formally decomposed as the sum of two time series, one deterministic and one stochastic.

$$x_t = \eta_t + \sum_{j=0}^{\infty} b_j \varepsilon_{t-j},\tag{5}$$

where  $\eta_t$  denotes the deterministic component and  $\varepsilon_t$  is the stochastic component that is input to an infinite vector of moving average weights  $\{b_j\}$ . Suppose Time series data in the same period can be regard as local weak-sense stationarity, The main benefit of weak-sense stationarity is that any time series can be put into the context Hilbert Space. It means any time series can be decomposed by a set of orthogonal increments in the space. Bochner's theorem (Loomis, 2013) ensures that there exists a group of Fourier-type complex exponential function  $\{e^{-2\pi i\xi t}\}$  componets to generate  $x_t$ .

$$x_t = \int e^{-2\pi i\xi t} d\omega_\xi \tag{6}$$

where  $d\omega_{\xi}$  is the measure weight with the  $\xi$ -th frequency wave. Due to the orthogonality of frequencies, it inspired us to explore whether the deterministic and stochastic components are dominated by specific frequencies, even in different periods and domains (detail study in Figure 3).

We first precompute FFT or DWT of each data in training set, calculate the averaged amplitude of each position in spectrum  $S = \{0, \dots, [L/2]\}$  and sort them by corresponding amplitude. Then we use  $\alpha$  to obtain top frequency positions  $S_a$  for decomposition in train and inference.  $S_\alpha$  can represent active stable frequency positions among different periods in multi-domain dataset, which has time-deterministic property (less affected by time) intuitively. Here we use  $\mathcal{F}$  to denote frequency transform operation,  $\mathcal{F}^{-1}$  denotes the inverse operation.

$$x_{sto} = \mathcal{F}^{-1} \left( \text{Mask} \left[ \mathcal{S}_{\alpha} \right] \cdot \mathcal{F}(x) \right), x_{det} = x - x_{sto}$$
<sup>(7)</sup>

where  $x_{det}$  denotes time-deterministic features and  $x_{sto}$  denotes time-stochastic features. Mask operation sets the amplitude of specific position components into zero.

### 247 5.2 TIME-DETERMINISTIC BLOCK

We use disentangled representation learning methods based on the paradigm of variational inference to deal with time-deterministic features  $x_{det}$ .

We construct the domain-invariant and domain-specific probabilistic encoders to decompose  $x_{det}$ feature space and a decoder merges the two parts of the features to reconstruct origin features. In parallel, it is supervised by constraint loss from both domain labels and category labels.

**Domain Disentangled Probabilistic Encoder.** Here we design domain-invariant and domainspecific probabilistic encoders denoted as  $q_{\phi}(z_{inv}|x_{det})$  and  $q_{\phi_d}(z_{spc}|x_{det})$  to represent the corresponding feature space seperately. The two encoders are used to learn the two feature space statistics ( $\mu_{inv}, \sigma_{inv}$ ) and ( $\mu_{spc}, \sigma_{spc}$ ). We use the reparametrization trick to resample  $z_{spc}$  and  $z_{inv}$ from standard normal distribution statistics.

259 260

261

262

$$q_{\phi}(z_{inv}|x_{det}) = \mathcal{N}\left(z_{inv} \mid \mu_{inv}, \sigma_{inv}; \phi\right), \tag{8}$$

$$q_{\phi_d}(z_{spc}|x_{det}) = \mathcal{N}\left(z_{spc} \mid \mu_{spc}, \sigma_{spc}; \phi_d\right),\tag{9}$$

where  $\phi$  and  $\phi_d$  denote the parameters of the domain-invariant encoder and domain-specific encoder respectively.

In order to adapt to prior distribution shift in different domains, we also set the prior hypothesis space  $p(z_{inv})$  and  $p(z_{spc})$  to data-driven Gaussian distributions by domain label and category label as well.

Finally, a probabilistic decoder  $p_{\theta}(x_{inv}|z_{inv}z_{spc})$  is used to feed  $z_{spc}$  and  $z_{inv}$  to reconstruct original information, where  $\theta$  are learnable parameters of the decoder. The probabilistic encoders and 270 decoder are learned with the following objective  $\mathcal{L}_{ELBO}$ : 271

273

274

275

285

286 287

288 289

291

292 293

295

296 297

298

299 300

301

$$\mathcal{L}_{ELBO} = \mathbb{E}_{q_{\phi}(z_{inv}|x_{det}), q_{\phi_d}(z_{spc}|x_{det})} \left[\log p_{\theta}(x_{inv}|z_{inv}, z_{spc})\right] - D_{\mathrm{KL}} \left(q_{\phi}(z_{inv} \mid x_{det}) \| p(z_{inv})\right) - D_{\mathrm{KL}} \left(q_{\phi_d}(z_{spc} \mid x_{det}) \| p(z_{spc})\right),$$
(10)

$$- D_{ ext{KL}} \left( q_{\phi_d}(z_{spc} \mid x_{det}) \| p(z_{spc}) 
ight)$$

where the first term represents the reconstruction loss to  $x_{det}$ , and the next two terms are the prior 276 matching terms from  $z_{inv}$  and  $z_{spc}$  by KL-divergence regularization. 277

278 Constraint Loss. The generation process under unsupervised signals proved to be unreliable (Lo-279 catello et al., 2019), so we use both domain labels and category labels for constraint loss. It encourages the domain-invariant space close to the label space, and the domain-specific space close to the domain space. We provide two schemes for constraint loss, one is the common feature cross-281 prediction loss (Qian et al., 2021), and the other is our proposed simplification of the constraint 282 process using supervised contrastive loss (Khosla et al., 2020), which achieves the same goal by 283 reducing the similarity of two features (Shi et al., 2024). 284

(1) feature cross-prediction loss.

$$\mathcal{L}_{con1} = \mathcal{L}_{li} + \mathcal{L}_{ds} - \mathcal{L}_{di} - \mathcal{L}_{ls}.$$
(11)

where  $\mathcal{L}_{li}$ ,  $\mathcal{L}_{ls}$  denotes  $z_{inv}$ ,  $z_{spc}$  with label loss,  $\mathcal{L}_{di}$ ,  $\mathcal{L}_{ds}$  denotes  $z_{inv}$ ,  $z_{spc}$  with domain class loss.

(2) domain-invirant contrastive loss. We use  $z_{inv}$  and  $z_{spc}$  of each batch of data itself as negative sample pairs, which not only increases the gap between the domain-invariant and domain-specific parts, but also reduces the individual differences between the positive samples.

$$\mathcal{L}_{con2} = -\sum_{i \in I} \log \frac{\exp(z_i, z^+/\tau)}{\sum_{k=0}^{K} \exp(z_i, z_k/\tau).}$$
(12)

where I is the origin index set of the  $z_{inv}$  and  $z_{spc}$  and we concat them to  $\{z_k\}_{k=0}^{K}$  features, and we made  $(z_i, \{z^+\}_{k=0}^{K^+}, \{z^-\}_{k=0}^{K^-})$ , while  $z^+$  represents the same feature space and the  $z^-$  represents the opposite space,  $\tau$  denotes to scalar temperature.

## 5.3 TIME-STOCHASTIC BLOCK

302 The time-stochastic block is constructed for time-stochastic features  $x_{sto}$ . It can be considered as 303 stochastic components with moving average weighting, so we try to capture the local time-variant 304 dynamic features with the help of patches (Nie et al., 2023), recent studies have shown that domain-305 invariant features in collaboration with other classification-related features help to improve the robustness of the model (Yu et al., 2024). 306

307 Patch Embedding. We divide the input time series into patches and set the patch length P and patch stride S, to divide the L-length sequence  $x_{sto}$  into M patches  $x_{patch} \in \mathbb{R}^{P \times M}$ , where  $M = \lfloor \frac{L-P}{S} \rfloor + 2$ . We map the patches through a linear layer  $W_p \in \mathbb{R}^{D \times P}$  to the transformer space with a learnable position encoding  $W_{pos} \in \mathbb{R}^{D \times M}$  to get the final patch embedding  $x_h \in \mathbb{R}^{D \times M}$ . 308 309 310 311

$$x_h = W_p x_{patch} + W_{pos}.$$
(13)

Stochastic Encoder. We use self-attention (Vaswani et al., 2017) in the transformer encoder to 314 model the dynamically varying weighting of the time dimension and to avoid mixing effects between 315 different variables, we follow a channel-independent design, taking each patch as input and final 316 splicing the outputs, and then an MLP classifier that constrains the classification loss in that part. 317 We keep stacking multi-head attention and feed-forward network with residual connections (He 318 et al., 2016) and layer normal (Ba et al., 2016).

$$Q_{h} = x_{h}^{T} W_{h}^{Q}, K_{h} = x_{h}^{T} W_{h}^{K}, V_{h} = x_{h}^{T} W_{h}^{V},$$
  

$$O_{h} = Attention(Q_{h}, K_{h}, V_{h}) = softmax(\frac{Q_{h} K_{h}^{T}}{\sqrt{d_{k}}})V_{h},$$
(14)

where  $O_h$  is the output from the multi-head attention layer and input to the feed-forward layer.

319 320

312

313

321

324 Finally, an MLP is utilized for label prediction, transferring the output of the encoder into the final 325 latent feature  $z_{std} \in \mathbb{R}^D$ , we use *cross-entrpy* (CE) loss  $L_{sto}$  for label classification. 326

327 328

329 330

331

333

335 336

337

338

342 343 344

345

347

348

349

350

351

352 353 354

360 361

362

363

364

365

366

#### MODEL SUMMARY AND THEORETICAL INSIGHTS 5.4

**Proposition 5.1.** Assuming time series data  $x_t \in \mathbb{R}^{L \times C}$  through frequency decomposition into time-deterministic components  $A_{det} = \{a_i\} \in \mathbb{R}^{K \times C}$  and time-stochastic components  $A_{sto} =$ 332  $\{a_i\} \in \mathbb{R}^{(L/2-K) \times C}$ , the invariant minimization objective formula suitable for OOD time series classification is as follows: 334

$$\min_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) := \sum_{e \in \mathcal{E}_{tr}} \mathcal{R}^{e}(\Phi_{det}(A_{det})) + \lambda_{det} \mathcal{P}(\Phi_{det}(A_{det})) + \lambda_{sto} \mathcal{J}_{\phi_{\boldsymbol{w}}}(A_{sto})$$
(15)

where  $\mathcal{R}^e$  denotes the risk in domain environment  $e, \Phi_{det}$  denotes invariant feature extractor.  $\mathcal{P}$  is the regularization for invariant feature,  $\mathcal{J}_{\phi_m}$  capture auxiliary features avoiding information lose.

339 FEDNet decoupled the feature into time-deterministic and time-stochastic features to solve domain 340 distribution shift and temporal distribution shift respectively and concat  $z_{inv}$ ,  $z_{sto}$  from two parts for final classification  $\mathcal{L}_{cls}$ . The model's total loss  $\mathcal{L}$  consists of three parts: 341

$$\mathcal{L} = \lambda_{det} \mathcal{L}_{det} + \lambda_{sto} \mathcal{L}_{sto} + \mathcal{L}_{cls},$$

$$\mathcal{L}_{det} = \mathcal{L}_{FLBO} + \mathcal{L}_{com},$$
(16)

where  $\mathcal{L}_{det}$  denotes time-deterministic loss,  $\mathcal{L}_{sto}$  denotes time-stochastic loss,  $\mathcal{L}_{cls}$  denotes final classification loss,  $\lambda_{det}$  and  $\lambda_{sto}$  are hyperparameter weights for two parts. 346

Proposition 5.2 (Frequency Perspective Risk Bound on Unseen Time series Domain). Let  $\mathcal{H}$  be a hypothesis space built from a set of source time series domains  $D = \{D_i\}_{i=1}^{N_d}$ . Suppose q > 0is a constant, for any unseen time series domain  $D_U$  from the convex hull  $\Lambda_D$ , we have its closest element  $D_{\bar{U}}$  related to source domains in  $\Lambda_D$ , i.e.,  $D_{\bar{U}} = \arg \min_{\pi_1, \dots, \pi_{N_d}} \beta_q(D_{\bar{U}} || \sum_{i=1}^{N_d} \pi_i D_i)$ . Then the risk of  $D_U$  on any label function  $h \in \mathcal{H}$  is,

$$R_{D_U}[h] \le \frac{1}{2} d_{D_U}(h) + \rho \cdot \left[ e_{D_{\bar{U}}}(h) \right]^{1 - \frac{1}{q}}, \tag{17}$$

where  $\rho = 2^{\frac{q-1}{q} \sup_{i,j \in [N_d]} RD_q(D_i || D_j)}$ ,  $d_D(h)$  and  $e_D(h)$  are ideal and empirical risk of domain D,

$$RD_{q}(D_{i}||D_{j}) = \frac{1}{q-1}\log\int \left[P_{F}(D_{i})\right]^{q} \left[P_{F}(D_{j})\right]^{1-q} da = \frac{1}{q-1}\log\left[\frac{\mu_{j}^{2(q-1)}}{\mu_{i}^{2q}} \cdot \frac{\sqrt{\pi}^{2n+1}n!!}{2^{2n+1}\sqrt{\gamma}^{n+1}}\right]$$
(18)

where  $\mu_i = E(D_i) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \sigma_k, \mu_j = E(D_j) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \tau_k, \gamma = \sum_{k=1}^n \frac{q}{2\sigma_k^2} + \frac{1-q}{2\tau_k^2}, \{\sigma_k\}_{k=1}^n$ and  $\{\tau_k\}_{k=1}^n$  denote frequency scale parameters in  $D_i, D_j, n$  represents the number of components.

The individual  $\sigma_k$  is a linear unbiased estimate of  $E(a_k)$ , The first term can be regarded as constant and bounded. The second term perfectly aligns with our motivation for decoupling the frequency domain. By keeping top  $E(a_k)$ , we can reduce n while making  $\gamma$  within a controllable range to decrease  $\rho$  and whole generalization bound. proofs for two proposition are provided in Appendix A

#### 6 EXPERIMENT

367 368

Datasets. We conduct serveral datasets to evaluate the performance and efficiency of FEDNet. We 369 used the synthetic dataset Spurious Fourier and the real datasets HHAR (Gagnon-Audet et al., 2022) 370 provided by WOODS (Gagnon-Audet et al., 2022), three open-source datasets UCIHAR (Anguita 371 et al., 2012), UniMiB-SHAR (Micucci et al., 2017), and Opportunity (Chavarriaga et al., 2013) 372 used in GILE(Qian et al., 2021), as well as DSADS (Barshan & Yüksek, 2014), PAMAP (Reiss & 373 Stricker, 2012), processed according to the domain division strategy provided by Diversify (Lu et al., 374 2022). The details of these datasets are listed in Appendix B. 375

Baselines. We evaluate the proposed FEDNet with various significant baselines, which can be 376 divided into three types. We provide detailed hyperparameter implementation, dataset settings and 377 comparing rules in Appendix D.

- General time series methods. We compare with the mainstream time series models in recent years such as PatchTST (Nie et al., 2023), and FreTS (Yi et al., 2023).
- General OOD methods. We choose some important baselines for OOD generalization from other research domains, GroupDRO (Sagawa\* et al., 2020), ANDMask (Parascandolo et al., 2021), and VREx (Krueger et al., 2021).
- OOD time series methods. Research on OOD generalization for time series classification is limited, and we select three important works, i.e., GILE (Qian et al., 2021), AdaRNN (Du et al., 2021), and Diversify (Lu et al., 2022).

6.1 PERFORMANCE COMPARISON

Table 1: Accuracy on cross-person generalization. "Target" represents the unseen test domain. Spurious Fourier is a synthetic dataset, with only  $\{d=10\%\}$  used as the test domain, while the remains are all real-world datasets.  $FEDNet_f$  uses FFT and  $FEDNet_w$  uses othogonal function DWT.

Dataset	Target	VREx	GroupDRO	ANDMasl	x FreTS	PatchTST	[GILE	AdaRNN	Diversify	FEDNet <sub>f</sub>	$FEDNet_w$
Spurious Fourier	d=10%	48.19	48.66	11.16	49.09	11.03	15.90	50.12	15.37	74.56	33.34
	0	89.60	88.49	91.48	93.12	79.46	94.51	72.77	91.48	96.97	<b>97.76</b>
ппур		91.00	88.81	92.65 86.81	93.76	79.90	96.94	75.67	92.65	98.26 00.76	97.27
IIIAK	3	65.55	65.26	61.60	62.76	39.68	63.34	50.06	54.32	63.17	67.87
	4	54.05	<u>55.74</u>	51.69	55.40	42.22	43.92	36.49	45.14	57.09	51.35
	0	89.34	89.34	98.56	81.55	78.96	83.07	80.20	87.03	94.81	78.38
	1	66.23	57.28	69.21	56.62	73.50	75.62	76.24	<u>76.49</u>	80.13	72.84
UCIHAR	2	97.65	96.77	<u>97.94</u>	92.30	70.67	86.19	86.45	90.91	97.95	90.61
	3	83.28	84.54	88.33	79.49	78.86	91.25	87.50	89.27	91.79	80.12
	4	/0.86	66.23	89.74	89.07	80.46	85.62	87.81	92.38	98.34	90.06
	1	55.99	<u>58.33</u>	57.55	41.14	49.21	47.39	46.88	50.26	55.98	64.58
UniMiB-SHAR	2	57.80	59.69	59.35	36.02	68.78	46.40	26.76	42.20	70.15	73.58
Children Strate	3	63.16	63.82	64.14	60.85	71.38	62.18	46.05	60.20	71.71	75.65
	5	41.28	42.62	40.94	38.25	36.91	38.43	35.57	44.30	40.26	<u>44.29</u>
	S1	53.73	59.62	77.85	81.87	52.00	84.02	80.64	82.23	84.86	83.82
Opportunity	S2	37.17	55.79	78.57	81.21	66.62	81.39	78.97	79.96	81.45	81.53
Opportunity	S3	36.77	56.31	74.67	75.94	46.76	<u>77.91</u>	76.36	76.79	79.11	76.32
	S4	46.41	58.15	78.32	77.58	52.48	80.91	78.85	80.74	81.77	80.60
	0	70.25	70.66	71.60	71.77	33.92	62.96	54.11	67.55	73.00	64.55
FMG	1	85.50	83.08	82.52	80.15	36.82	68.02	57.44	81.09	87.10	59.59
LIVIO	2	73.62	77.03	76.91	74.88	22.66	66.02	57.83	74.64	79.66	77.51
	3	77.14	78.62	77.50	77.96	36.62	69.99	53.87	77.32	77.43	79.85
	0	80.26	84.69	82.50	80.26	82.24	89.64	83.11	77.19	92.80	92.41
DEADS	1	76.54	78.03	73.42	70.13	74.07	78.20	79.78	77.28	84.86	83.64
DOADO	2	86.40	85.96	83.03	84.29	82.67	86.75	83.46	85.22	93.24	90.65
	3	74.61	74.39	78.46	73.46	78.85	79.56	70.35	71.80	87.71	80.52
	0	62.88	61.75	61.84	55.22	60.40	65.01	63.30	61.98	64.94	67.48
DAMAD	1	54.88	52.00	53.04	60.41	66.36	51.46	54.24	54.38	68.16	67.08
PAMAP	2	22.68	25.69	28.02	34.98	50.06	25.23	23.35	24.32	34.39	35.27
	3	62.10	65.22	67.86	68.69	63.39	68.06	61.04	57.79	67.13	<u>68.55</u>

421

378

379

380

381

382

384

386 387

388 389

390

391

392

422 We evaluate the generalizability of our proposed FEDNet by comparing its performance with base-423 lines on these publicly accessible datasets in Table 16. We choose accuracy as the main evaluation metric. For our experiments, we select one domain as the "target domain" for testing, while the 424 remaining domains serve as the "source domains" for training. 425

426 **Cross-person generalization.** The frequency-domain separation design of FEDNet achieves the 427 best performance across multiple datasets compared to other methods. General time series methods 428 use ERM (Empirical Risk Minimization) as the optimization objective, resulting in greater perfor-429 mance fluctuations across different domains. These methods are significantly affected by domain shifts and segmentation. FreTS, which also extracts frequency-domain information, performs bet-430 ter than PatchTST in some domains, indicating that frequency information is robust, by providing 431 a global perspective on time series. Additionally, FEDNet outperforms three different representative types of General OOD that fail to consider the impact of temporal distribution shifts within
the time series window on marginal probabilities. Methods like AdaRNN and Diversify, which are
carefully designed to account for temporal shifts through reweighting or relabeling, directly extract
information from the pure time domain. Essentially, these are data augmentation techniques and are
inevitably affected by sample noise when achieving domain generalization. We also provide results
for other metrics, as detailed in the Appendix C.9.

438
 439
 439
 440
 440
 440
 441
 441
 441
 442
 442
 442
 442
 443
 444
 444
 444
 444
 445
 446
 447
 447
 448
 448
 449
 449
 440
 440
 441
 441
 441
 441
 442
 442
 442
 442
 443
 444
 444
 444
 444
 444
 445
 446
 446
 447
 447
 447
 448
 448
 449
 440
 441
 441
 441
 441
 441
 442
 442
 442
 441
 441
 442
 442
 441
 441
 442
 442
 442
 442
 442
 442
 442
 441
 442
 442
 441
 442
 442
 441
 442
 441
 442
 442
 441
 442
 442
 442
 442
 442
 442
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444
 444

Table 2: Accuracy on cross-position generalization for DSADS.

Dataset	Target	VREx	GroupDRO	ANDMask	IB-IRM	IRM	GILE	AdaRNN	Diversify	$FEDNet_f$
DSADS	$\begin{array}{c}0\\1\\2\\3\\4\end{array}$	27.12 20.33 27.17 24.78 17.24	29.99 23.86 38.20 24.00 26.01	26.64 24.62 33.22 26.95 20.18	29.82 23.17 35.72 19.96 22.15	25.30 20.18 25.62 27.34 18.25	38.8521.4538.3320.0923.14	38.19 29.04 32.72 24.61 19.50	<b>47.70</b> <u>32.90</u> <b>44.50</b> <u>31.60</u> <u>30.40</u>	40.54 36.05 36.25 33.93 33.04

## 6.2 MODEL ANALYSIS

**Frequency Decomposition Study.** We add a fully connected linear layer behind two components separately after the Frequency Filter to analyze the linear weight changes. We chose different domains as the target domains for training and computed the coefficient of variation of the weights of the two linear layers since the linear weights can reflect the dependence between the time series points. Formally, the coefficient of variation is defined as follows:

coefficient of variation 
$$= \frac{\sigma_{weights}}{\mu_{weights}},$$
 (19)

where  $\sigma_{weight}$  and  $\mu_{weights}$  represent the standard deviation and mean of the linear layer weights.

As shown in Figure 3, we find that the coefficient of variation of the deterministic linear weights during training in different domains is always lower than that in stochastic part. It shows that the deterministic component could vary less among different periods, controlled by specific frequency.





Invirant Feature Study. We use A-distance (Schölkopf et al., 2007) to measure the domain discrepancy of invariant fea-tures obtained by different methods. It can be approximated as  $d = 2(1 - 2\sigma_A)$ , where  $\sigma_A$  is the risk of a binary classi-fier distinguishing features between source and target domains. Figure 4 shows that our method consistently outperforms other approaches, and using contrastive learning achieves better re-sults compared to the standard cross label loss. The smaller the indicator, the more invariant features are.



Figure 4: A-distance on invariant features with EMG dataset.

**Ablation Study.** we conduct ablations on FEDNet with several degenerate variants to analyze model components: (1) w/o  $\mathcal{L}_{con}$ : we remove the constraint loss in Time-Determistic Block. (2) w/o  $\mathcal{L}_{det}$ : we remove the time-determistic block, only use  $x_{sto}$  for classification. (3) w/o  $\mathcal{L}_{sto}$ : we remove the time-stochastic block, only use  $x_{det}$  for classification. We present some of the results in Figure 5, with detailed results available in Appendix C.1. We observe that our proposed Frequency Filter is effective in disentangling distributional shifts, and it obtains stable time-invariant components that eliminate the effect of temporal shifts:

493

508

509

510

511

512

518 519 520

521

522 523

524

525

526

527

528

494 (1) w/o  $\mathcal{L}_{con}$ : The time-deterministic block es-495 sentially degenerates into a dimensionality re-496 duction module for  $x_{det}$  without considering 497 domain distribution shift. The final results are 498 not much different from the general time series 499 methods.

500 (2) w/o  $\mathcal{L}_{det}$ : Lack of  $x_{det}$  leads to significant 501 performance degradation. At the same time, 502 we find that time-stochastic features still have 503 some classification potential in some domains.

(3) w/o  $\mathcal{L}_{sto}$ : without  $x_{sto}$ , it also has a small effect on model performance. It shows that dy-



Figure 5: Ablation study. The Y-axis shows average accuracy.

namic changes in the time series can improve the robustness of the model in unseen domains tosome degree.

(4) cross vs contrast: contrastive learning obviously outperforms cross loss.

**Empirical Domain Divergence Study.** As Proposition 5.2 shows,  $\gamma$  can be seen as a meaningful metric to measure the volatility of time series in frequency. We adopted it to be an empirical absolute value  $\hat{\gamma} = \sum_{k=1}^{n} q \left| \frac{1}{2\sigma_k^2} - \frac{1}{2\tau_k^2} \right|$  avoid redundant frequency to estimate the domain divergence. Figure 3 (d) changes of domain divergence in the dataset consistents with the IRM gains experiments.



Figure 6: (a) denotes the second log value with the number of frequencies n and  $\gamma$ . It shows that we could get the lowest upper bound on generalization when  $\alpha \in [0.2, 0.4]$ ,  $\gamma \in [1, 20]$  is usually the datasets' range after reducing frequencies. (b) and (c) denote the trend of  $\hat{\gamma}$  when q = 2, 1/2, shows that reducing the number of frequencies reduces the volatility of time series. (d) denotes dataset domain divergence with mask  $\alpha$ =0.2. It demonstrates that keep  $\alpha$  ratio high amplitude frequency reduces domain divergence. The anomaly in SHAR stems from denominator  $\hat{\gamma}(< 0.01)$  very small.

529 530 531

532

## 7 CONCLUSION

In this paper, we focus on OOD time series classification and propose a novel method called FED-Net. Our method incorporates frequency information as a prior and utilizes a decomposition frame-work to separate time series into time-deterministic components and time-stochastic components. We address both domain distribution shift and temporal distribution shift to extract invariant features for domain generalization. Extensive experiments demonstrate that FEDNet achieves superior performance and effectively exploits frequency information for OOD time series classification.

## 540 REFERENCES

- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio,
  Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-ofdistribution generalization. *Advances in Neural Information Processing Systems*, 34:3438–3450,
  2021.
- Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge L Reyes-Ortiz. Human activity recognition on smartphones using a multiclass hardware-friendly support vector machine. In *Ambient Assisted Living and Home Care: 4th International Workshop, IWAAL 2012, Vitoria-Gasteiz, Spain, December 3-5, 2012. Proceedings 4*, pp. 216–223. Springer, 2012.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.
   *arXiv preprint arXiv:1907.02893*, 2019.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Billur Barshan and Murat Cihan Yüksek. Recognizing daily and sports activities in two open source machine learning environments using body-worn sensor units. *The Computer Journal*, 57(11): 1649–1667, 2014.
- Ricardo Chavarriaga, Hesam Sagha, Alberto Calatroni, Sundara Tejaswi Digumarti, Gerhard Tröster, José del R Millán, and Daniel Roggen. The opportunity challenge: A benchmark database for on-body sensor-based activity recognition. *Pattern Recognition Letters*, 34(15):2033–2042, 2013.
- Berken Utku Demirel and Christian Holz. Finding order in chaos: A novel data augmentation method for time series in contrastive learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Don Dennis, Durmus Alp Emre Acar, Vikram Mandikal, Vinu Sankar Sadasivan, Venkatesh Saligrama, Harsha Vardhan Simhadri, and Prateek Jain. Shallow rnn: accurate time-series classification on resource constrained devices. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuntao Du, Jindong Wang, Wenjie Feng, Sinno Pan, Tao Qin, Renjun Xu, and Chongjun Wang.
  Adarnn: Adaptive learning and forecasting for time series. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM)*, 2021.
- Jean-Christophe Gagnon-Audet, Kartik Ahuja, Mohammad-Javad Darvishi-Bayazi, Pooneh
   Mousavi, Guillaume Dumas, and Irina Rish. Woods: Benchmarks for out-of-distribution gen eralization in time series. *arXiv preprint arXiv:2203.09978*, 2022.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
  Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016.
- Pascal Germain, Amaury Habrard, François Laviolette, and Emilie Morvant. A new pac-bayesian perspective on domain adaptation. In *International conference on machine learning*, pp. 859–868.
  PMLR, 2016.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G
   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 signals. *circulation*, 101(23):e215–e220, 2000.
- Huan He, Owen Queen, Teddy Koker, Consuelo Cuevas, Theodoros Tsiligkaridis, and Marinka Zit nik. Domain adaptation for time series under feature and label shifts. In *International Conference* on Machine Learning, pp. 12746–12774. PMLR, 2023.

605

609

615

622

628

630

594	Kaiming He. Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
595	nition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp.
596	770–778, 2016.
597	

- Maximilian Ilse, Jakub M Tomczak, Christos Louizos, and Max Welling. Diva: Domain invariant 598 variational autoencoders. In Medical Imaging with Deep Learning, pp. 322–348. PMLR, 2020.
- 600 Gwilym M. Jenkins, Herman O. A. Wold, and Peter Whittle. A study in the analysis of sta-601 tionary time-series. 1955. URL https://api.semanticscholar.org/CorpusID: 602 65334831. 603
  - Johan Ludwig William Valdemar Jensen. Sur les fonctions convexes et les inégalités entre les valeurs moyennes. Acta mathematica, 30(1):175-193, 1906.
- Lars Kegel, Martin Hahmann, and Wolfgang Lehner. Feature-based comparison and generation 607 of time series. In Proceedings of the 30th international conference on scientific and statistical 608 database management, pp. 1-12, 2018.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron 610 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. Advances in neural 611 information processing systems, 33:18661–18673, 2020. 612
- 613 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint 614 arXiv:1312.6114, 2013.
- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai 616 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapo-617 lation (rex). In International Conference on Machine Learning, pp. 5815–5826. PMLR, 2021. 618
- 619 Daniel Kuhn, Peyman Mohajerin Esfahani, Viet Anh Nguyen, and Soroosh Shafieezadeh-Abadeh. 620 Wasserstein distributionally robust optimization: Theory and applications in machine learning. In 621 Operations research & management science in the age of analytics, pp. 130–166. Informs, 2019.
- Bo Li, Yifei Shen, Yezhen Wang, Wenzhen Zhu, Dongsheng Li, Kurt Keutzer, and Han Zhao. In-623 variant information bottleneck for domain generalization. In Proceedings of the AAAI Conference 624 on Artificial Intelligence, volume 36, pp. 7399-7407, 2022. 625
- 626 Yong Liu, Chenyu Li, Jianmin Wang, and Mingsheng Long. Koopa: Learning non-stationary time 627 series dynamics with koopman predictors. *arXiv preprint arXiv:2305.18803*, 2023.
- Sergey Lobov, Nadia Krilova, Innokentiy Kastalskiy, Victor Kazantsev, and Valeri A Makarov. La-629 tent factors limiting the performance of semg-interfaces. Sensors, 18(4):1122, 2018.
- 631 Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard 632 Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning 633 of disentangled representations. In *international conference on machine learning*, pp. 4114–4124. 634 PMLR, 2019.
- Lynn H Loomis. Introduction to abstract harmonic analysis. Courier Corporation, 2013. 636
- 637 Wang Lu, Jindong Wang, Xinwei Sun, Yiqiang Chen, and Xing Xie. Out-of-distribution representa-638 tion learning for time series classification. In The Eleventh International Conference on Learning 639 Representations, 2022. 640
- Daniela Micucci, Marco Mobilio, and Paolo Napoletano. Unimib shar: A dataset for human activity 641 recognition using acceleration data from smartphones. Applied Sciences, 7(10):1101, 2017. 642
- 643 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 644 64 words: Long-term forecasting with transformers. In International Conference on Learning 645 Representations, 2023. 646
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. IEEE Transactions on knowledge 647 and data engineering, 22(10):1345-1359, 2009.

648 649 650	Giambattista Parascandolo, Alexander Neitz, ANTONIO ORVIETO, Luigi Gresele, and Bernhard Schölkopf. Learning explanations that are hard to vary. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=hblsDDSLbV.
651 652 653 654	Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain adaptation: A survey of recent advances. <i>IEEE Signal Processing Magazine</i> , 32(3):53–69, 2015. doi: 10. 1109/MSP.2014.2347059.
655 656 657	Hangwei Qian, Sinno Jialin Pan, and Chunyan Miao. Latent independent excitation for generaliz- able sensor-based cross-person activity recognition. In <i>Proceedings of the AAAI Conference on</i> <i>Artificial Intelligence</i> , volume 35, pp. 11921–11929, 2021.
658 659 660	Attila Reiss and Didier Stricker. Introducing a new benchmarked dataset for activity monitoring. In 2012 16th international symposium on wearable computers, pp. 108–109. IEEE, 2012.
661 662 663	Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In <i>International Conference on Learning Representations</i> , 2020. URL https://openreview.net/forum?id=ryxGuJrFvS.
664 665 666 667	Sohrab Salehi, Farhia Kabeer, Nicholas Ceglia, Mirela Andronescu, Marc J Williams, Kieran R Campbell, Tehmina Masud, Beixi Wang, Justina Biele, Jazmine Brimhall, et al. Clonal fitness inferred from time-series modelling of single-cell cancer genomes. <i>Nature</i> , 595(7868):585–590, 2021.
669 670	Bernhard Schölkopf, John Platt, and Thomas Hofmann. Analysis of Representations for Domain Adaptation, pp. 137–144. 2007.
671 672 673 674 675	Ruize Shi, Hong Huang, Kehan Yin, Wei Zhou, and Hai Jin. Orthogonality matters: Invariant time series representation for out-of-distribution classification. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , KDD '24, pp. 2674–2685, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704901. doi: 10. 1145/3637528.3671768. URL https://doi.org/10.1145/3637528.3671768.
676 677 678	Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log- likelihood function. <i>Journal of statistical planning and inference</i> , 90(2):227–244, 2000.
679 680 681 682	Akara Supratak, Hao Dong, Chao Wu, and Yike Guo. Deepsleepnet: a model for automatic sleep stage scoring based on raw single-channel eeg. <i>IEEE Transactions on Neural Systems and Rehabilitation Engineering</i> , 25(11):1998–2008, Nov 2017. ISSN 1534-4320. doi: 10.1109/TNSRE.2017.2721116.
683 684 685	Chi Ian Tang, Ignacio Perez-Pozuelo, Dimitris Spathis, and Cecilia Mascolo. Exploring contrastive learning in human activity recognition for healthcare. <i>arXiv preprint arXiv:2011.11542</i> , 2020.
686 687 688	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.
689 690 691 692	Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and S Yu Philip. Generalizing to unseen domains: A survey on domain generalization. <i>IEEE transactions on knowledge and data engineering</i> , 35(8):8052–8072, 2022a.
693 694 695	Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. Generalizing to unseen domains: A survey on domain generalization. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2022b.
696 697 698	Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. A fourier-based framework for domain generalization. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 14383–14392, 2021.
700 701	Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Longbing Cao, and Zhendong Niu. Frequency-domain MLPs are more effective learners in time series forecasting. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023.

702 703 704	Xi Yu, Huan-Hsin Tseng, Shinjae Yoo, Haibin Ling, and Yuewei Lin. Insure: an information theory inspired disentanglement and purification model for domain generalization. <i>IEEE Transactions on Image Processing</i> , 2024.
705 706 707	Hongyi Zhang. mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412, 2017.
708 709 710	Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. Self-supervised con- trastive pre-training for time series via time-frequency consistency. In <i>Proceedings of Neural</i> <i>Information Processing Systems, NeurIPS</i> , 2022.
711 712 713 714	Xingxuan Zhang, Peng Cui, Renzhe Xu, Linjun Zhou, Yue He, and Zheyan Shen. Deep stable learning for out-of-distribution generalization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 5372–5382, 2021.
715 716 717	Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In <i>Proc. 39th International Conference on Machine Learning (ICML 2022)</i> , 2022.
718 719 720	
721 722 723	
724 725 726	
727 728 729	
730 731 732	
733 734	
735 736 737	
738 739 740	
741 742 743	
744 745 746	
747 748 749	
750 751	
752 753 754 755	
100	

## A THEORETICAL INSIGHTS

# 758 A.1 BACKGROUND

Here, we first provide some background knowledge to illustrate the widespread occurrence ofmarginal probability shift in time series.

**OOD Generalization Problem.** Given datasets  $D_e = \{(x_i^e, y_i^e)\}_{i=1}^{n_e}$  collected from multiple training environments  $e \in \mathcal{E}_{tr}$ . Each dataset  $D_e$  contains a group of examples according to a certain probability distribution  $x \sim \mathbb{P}^e(x)$ . We hold on the basic covariate shift assumption (Shimodaira, 2000) that the optimal label function  $h^* \in \mathcal{H}$  representation space or conditional probability distribution  $\mathbb{P}^e(y|x)$  are the same within different environments while  $\mathbb{P}^e(x) \neq \mathbb{P}^{e'}(x), \forall e, e' \in \mathcal{E}_{tr}$ .

We define the prediction risk of model f under e environment as  $R^e(f) := \mathbb{E}^e \left[ \ell \left( f \left( x^e \right), y^e \right) \right]$ . Our goal is to obtain a optimal model  $f^*_{\theta}$  that can generalize to unseen domain distribution  $\mathcal{E}_{unseen} = \mathcal{E}_{all} / \mathcal{E}_{tr}$ .

767

768

769

 $f_{\theta}^* = \arg\min_{f_{\theta}} \mathbb{E}_{\mathcal{E}_{unseen}}[\ell(f_{\theta}(x), y)].$ (20)

789

790

791

792

793 794

796 797

798

799

800

806

808

where  $\mathbb{E}$  denotes expectation and  $\ell(\cdot, \cdot)$  denotes loss function.

775 **OOD Time Series Problem.** The basic goal of the Time Series Out of Distribution Problem is the 776 same as that of the Out of Distribution Generalization Problem. However, the major difference is 777 that there exists temporal distribution shift in each domain, which causes the data distribution broken 778 the i.i.d assumption with  $D_e(x, y)$ , i.e.  $\mathbb{P}^{D_k^i}(x, y) \neq \mathbb{P}^{D_k}(x, y), \exists i \neq j \in [1, T]$ .

**IRM.** Existing IRM methods encourage to elicit an invariant predictor f(:, w) which is a composite function of  $w \circ \Phi$  across environments  $\mathcal{E}_{all}$ . The feature extractor  $\Phi : \mathcal{X} \to \mathcal{H}$  maps  $\mathcal{X}$  to representation space  $\mathcal{H}$  to extract invariant features h from  $\mathcal{E}_{tr}$  which support  $\mathbb{E}[y^e \mid \Phi(x^e) = h] = \mathbb{E}\left[y^{e'} \mid \Phi(x^{e'}) = h\right], \forall e, e' \in \mathcal{E}_{all}$ , while the classifier  $w : \mathcal{H} \to \mathcal{Y}$  simultaneously optimal the prediction in among  $\mathcal{E}_{tr}$ . Here is the formula minimization objective:

$$\min_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) := \sum_{e \in \mathcal{E}_{tr}} \mathcal{R}^{e}(\boldsymbol{w}) + \lambda \mathcal{P}(\boldsymbol{w})$$
(21)

where  $\mathcal{R}^{e}(w) = \frac{1}{n_{e}} \sum_{i=1}^{n_{e}} \ell(f(\mathbf{x}_{i}^{e}, w), \mathbf{y}_{i}^{e})$  and  $\ell$  is loss function.  $\mathcal{P}(w)$  is a regularization encourage f(:, w) to optimal all environments. IRM is mainly used to solve the problem of distributional shifts due to conditional probabilities  $\mathbb{P}^{e}(y|x)$  in different environments, hoping to find stable and invirant features in different environments to solve the OOD problem, but IRM will face challenges with huge marginal probability shifts  $\mathbb{P}^{e}(x)$ .

## A.2 PRELIMINARY

**Lemma 2 (Temporal Covariate Shift).** The root cause of Temporal Distribution Shift is the  $\mathbb{P}(x)$  marginal distribution changes, while the  $\mathbb{P}(y|x)$  conditional distribution in each domain remains unchanged.

Lemma 3 (Phase Congruency). In Frequency distribution shift, the change in the phase is small and can be ignored.

Lemma 4 (Distribution of Fourier Component). The distributions of Fourier amplitude and phase
 can be modeled as Rayleigh distribution and uniform distribution respectively (He et al., 2023), the
 probabilistic density function of fourier component can be formulated:

$$f(a,p) = \operatorname{Rayleigh}(a|\sigma) \cdot \operatorname{U}(p|0,2\pi) = \frac{a}{2\pi\sigma^2} \cdot \exp(-\frac{a^2}{2\sigma^2}),$$
(22)

where a denotes the amplitude while p represents phase,  $\sigma$  is a variance parameter to scale the distribution of frequency component.

The marginal probabilistic density function of one fourier component amplitude can be viewed as the integral of f(a, p) over  $p \in [0, 2\pi], a \in (0, +\infty)$ 

$$f(a) = \int_{0}^{2\pi} f(a, p) dp = \frac{a}{\sigma^2} \cdot \exp(-\frac{a^2}{2\sigma^2}) da$$
(23)

This also confirms the validity of the influence of linear phase changes on the probability distribution as stated in **Lemma 3**.

Then we can get the expectation E(a) and variance D(a) of the frequency components

$$E(a) = \sqrt{\frac{\pi}{2}}\sigma, D(a) = \frac{4-\pi}{2}\sigma^2.$$
(24)

We can easily calculate the conclusion that E(a) is an unbiased estimate of the statistic, which means that we can directly reflect the change on the  $\sigma$  parameter by the average amplitude of frequency component, It is consistent with our motivation for magnitude decoupling.

**Corollary 1 (Frequency Marginal Probability Distribution of Time Series Data).** Combining **Lemma 3** and **Lemma 5** with frequency components independent assumption(Kegel et al., 2018), we can obtain the strict formulation of the frequency marginal probability distribution  $\mathbb{P}_F(x)$  given a time series data in any domain,

$$\mathbb{P}_{F}(x) = \int f(F[1], \dots F[n]) dF = \int f(a_{1}, \dots, a_{n}) = \int \prod_{k=1}^{n} f(a_{k}) da = \int \prod_{k=1}^{n} \frac{a}{\sigma_{k}^{2}} \cdot \exp(-\frac{a^{2}}{2\sigma_{k}^{2}}) da,$$
(25)

Naturally, we can use the principle of probabilistic independence to obtain the overall expectation and variance of the time series in the frequency domain perspective.

$$E(x) = \prod_{k=1}^{n} E(a_k) = \left(\frac{\pi}{2}\right)^{\frac{n}{2}} \prod_{i=k}^{n} \sigma_k, D(x) = \sum_{k=1}^{n} D(a_k) = \frac{4-\pi}{2} \sum_{k=1}^{n} \sigma_k^2,$$
 (26)

**Lemma 5 (Domain Divergence).(Germain et al., 2016) Suppose any domain**  $D_1$ ,  $D_2$  are built on input variable x and label variable y. Let q > 0 be a constant, the domain divergence between  $D_1$  and  $D_2$  is defined as

$$\beta_q \left( D_1 \| D_2 \right) = \left[ \mathbb{E}_{(x,y) \sim D_2} \left( \frac{D_1(x,y)}{D_2(x,y)} \right)^q \right]^{\frac{1}{q}} = 2^{\frac{q-1}{q}RD_q(D_1\|D_2)}, \tag{27}$$

where  $RD_q(\cdot)$  is Rényi Divergence.

**Corollary 2 (Bounding Domain Divergence in A Convex Hull).** Let D be a set of source domains, denoted as  $D = \{D_i\}_{i=1}^{N_d}$ . A convex hull  $\Lambda_D$  considered here consists of a mixture of distributions  $\Lambda_D = \left\{ \overline{D} : \overline{D}(\cdot) = \sum_{i=1}^{N_d} \pi_i D_i(\cdot), \ \pi_i \in \Delta_{N_d} \right\}$ , where  $\Delta_{N_d}$  is the  $N_d$  – 1-th dimensional simplex  $(\forall \pi_i > 0, \sum_i^{N_d} \pi_i = 1)$ . Let  $\rho = 2^{\frac{q-1}{q} \sup_{i,j \in [N_d]} RD_q(D_i || D_j)}$  then we have the following relation for the domain divergence when q > 1 between any pair of two domains  $D', D'' \in \Lambda_S$  in the convex hull,

$$\beta_q(D' \parallel D'') \le \rho. \tag{28}$$

*Proof.* Suppose two unseen domains D' and D'' on the convex hull  $\Lambda_S$  of  $N_d$  source domains with support  $\mathcal{O}$ . More specifically, let these two domains be  $D' = \sum_{i=1}^{N_d} \pi_i D_i(\cdot)$  and  $D'' = \sum_{i=1}^{N_d} \pi_j D_j(\cdot)$ , then the domain divergence between D' and D'' is

$$\beta_q \left( D' \| D'' \right) = 2^{\frac{q-1}{q} R D_q \left( D' \| D'' \right)}. \tag{29}$$

Let us consider the part of  $RD_q(\cdot)$  as follows first,

$$RD_q(D'||D'') = \frac{1}{q-1} \ln \int_{\mathcal{O}} \left[ \sum_{i=1}^{N_d} \pi_i D_i(x) \right]^q \left[ \sum_{j=1}^{N_d} \pi_j D_j(x) \right]^{1-q} dx$$
(30)

Since  $\forall \pi_i \in \Delta_{N_d}$ , We knows that  $\sum_{i=1}^{N_d} \pi_i = 1, \forall \pi_i \ge 0$ , when q > 1 we could hold  $f(x) = x^q$ and  $g(x) = x^{1-q}$  are convex function, their second-order derivative coefficients are the same and positive  $q \times (q-1) > 0$ .

Thus the original equation satisfies Jensen's Inequality (Jensen, 1906) and we get the following inequality

$$RD_{q}(D'\|D'') = \frac{1}{q-1} \ln \int_{\mathcal{O}} \left[ \sum_{i=1}^{N_{d}} \pi_{i} D_{i}(x) \right]^{q} \left[ \sum_{j=1}^{N_{d}} \pi_{j} D_{j}(x) \right]^{1-q} dx$$
(31)

$$\leq \frac{1}{q-1} \ln \int_{\mathcal{O}} \sum_{i=1}^{N_d} \pi_i \left[ D_i(x) \right]^q \sum_{i=j}^{N_d} \pi_j \left[ D_j(x) \right]^{1-q} dx \tag{32}$$

$$\leq \frac{1}{q-1} \ln \sum_{i=1}^{N_d} \sum_{i=j}^{N_d} \pi_i \pi_j \int_{\mathcal{O}} \left[ D_i(x) \right]^q \left[ D_j(x) \right]^{1-q} dx \tag{33}$$

$$\leq \frac{1}{q-1} \ln \int_{\mathcal{O}} \left[ D_i(x) \right]^q \left[ D_j(x) \right]^{1-q} dx \tag{34}$$

$$\leq \sup_{i,j \in [N_d]} RD_q(D_i \| D_j) \tag{35}$$

Then we have

$$\beta_q \left( D' \| D'' \right) = 2^{\frac{q-1}{q}RD_q(D'\|D'')} \le 2^{\frac{q-1}{q}\sup_{i,j\in[N_d]}RD_q(D_i\|D_j)} = \rho.$$
(36)

## A.3 PROOF OF PROPOSITION 5.1

**Proposition 5.1.** Assuming time series data  $x_t \in \mathbb{R}^{L \times C}$  through frequency decomposition into time-deterministic components  $A_{det} = \{a_i\} \in \mathbb{R}^{K \times C}$  and time-stochastic components  $A_{sto} = \{a_i\} \in \mathbb{R}^{(L/2-K) \times C}$ , the invariant minimization objective formula suitable for OOD time series classification is as follows:

$$\min_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) := \sum_{e \in \mathcal{E}_{tr}} \mathcal{R}^{e}(\Phi_{det}(A_{det})) + \lambda_{det} \mathcal{P}(\Phi_{det}(A_{det})) + \lambda_{sto} \mathcal{J}(A_{sto}),$$
(37)

where  $\mathcal{R}^e$  denotes the risk in domain environment  $e, \Phi_{det}$  denotes invariant feature extractor.  $\mathcal{P}$  is the regularization for invariant feature,  $\mathcal{J}_{\phi_w}$  capture auxiliary features avoiding information lose.

*Proof.* Since the frequency components obtained through the Fourier transform are orthogonal, we can consider these frequencies to be probabilistically independent of each other.

$$F[k] = \mathcal{F}[f(t)] = \int_{-\infty}^{\infty} f(t)e^{-i2\pi kt}dt = A_k e^{jP_k}$$
(38)

where  $A_k$  is the amplitude and  $P_k$  is the phase of the k-th frequency conponent, thus we obtain:

$$\forall i, j, F[i] \perp F[j] \to \mathbb{P}_F(F[i]) \perp \mathbb{P}_F(F[j])$$
(39)

From Lemma 2, Lemma 3, the phase is not influenced by distributional shifts, so the probability of frequency components is determined by the amplitude frequency.

$$\mathbb{P}_F(x) = \prod_{k=1}^K \mathbb{P}_F(F[k]) = \prod_{k=1}^K \mathbb{P}_F(A_k)$$
(40)

915 After Frequency Filter, the time series decomposed into the time-deterministic components  $A_{det} = \{a_1, a_2, \dots, a_k\}$  and time-stochastic components  $A_{sto} = \{a_{k+1}, \dots, a_{L/2}\}$ . Then we could rewrite 917 the  $\mathbb{P}_F(x)$ 

$$\mathbb{P}_F(x) = \mathbb{P}_F(A_{det}) \times \mathbb{P}_F(A_{sto}) \tag{41}$$

The probability of  $\mathbb{P}_F(A_{det})$  can be considered unaffected by temporal influences cross domains, i.e  $\forall 1 \leq i \neq j \leq T$ , SUPP $(\mathbb{P}_F^{D_k^i}(A_{det})) =$  SUPP $(\mathbb{P}_F^{D_k^j}(A_{det}))$ , thus making this feature suitable for the IRM invariance theory. thus making this feature suitable for the IRM invariance theory by holding on the conidtion assumption, 

$$\mathbb{E}[y^e \mid \Phi(A_{det})] = \mathbb{E}[y^{e'} \mid \Phi(A_{det})], \forall e, e' \in \mathcal{E}_{all}$$
(42)

At the same time, despite  $\forall 1 \leq i \neq j \leq T$ , SUPP $(\mathbb{P}_F^{D_k^i}(A_{sto})) \neq$  SUPP $(\mathbb{P}_F^{D_k^j}(A_{sto}))$  could lead to the frequency distribution shift, the behavior on the time series is influenced by moving average weights, suppose a learnable feature reweight extractor could alleviate the shift optimized with ERM  $J_{\phi_m}(\cdot)$  since  $\mathbb{P}^{D_k}(y|x)$  don't change ensure its practicability to help generalization space without affecting  $\Phi_{det}$  extractor and losing information.

$$\min_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) := \sum_{e \in \mathcal{E}_{tr}} \mathcal{R}^{e}(\Phi_{det}(A_{det})) + \lambda_{det} \mathcal{P}(\Phi_{det}(A_{det})) + \lambda_{sto} \mathcal{J}_{\phi_{\boldsymbol{w}}}(A_{sto})$$
(43)

Therefore, we complete the proof.

#### THE RISK FOR DOMAIN GENERALIZATION OF TIME SERIES FROM FREQUENCY VIEW A.4

Theorem 1 (PAC-Bayesian Risk Bound on Unseen Time series Domain).(Germain et al., 2016) Let  $\mathcal{H}$  be a hypothesis space built from a set of source domains, denoted as  $D = \{D_i\}_{i=1}^{N_d}$ . Suppose q > 0 is a constant, for any unseen domain  $D_U$  from the convex hull  $\Lambda_D$ , we have its closest element  $D_{\bar{U}}$  in  $\Lambda_D$ , i.e.,  $D_{\bar{U}} = \arg \min_{\pi_1, \dots, \pi_{N_d}} \beta_q(D_{\bar{U}} || \sum_{i=1}^{N_d} \pi_i D_i)$ . Then the risk of  $D_U$  on any label function  $h \in \mathcal{H}$  is,

$$R_{D_U}[h] \le \frac{1}{2} d_{D_U}(h) + \epsilon \cdot \left[ e_{D_{\bar{U}}}(h) \right]^{1 - \frac{1}{q}} + \eta_{T/S},\tag{44}$$

where  $d_D(h)$  and  $e_D(h)$  are ideal and expected risk of a domain D respectively,  $\epsilon =$  $\beta_q(D_U \| \sum_{i=1}^{N_d} \pi_i D_i)$  is an ideal distance since we can't have access to  $D_U$ , while  $\eta_{T/S}$  denotes the distribution of  $(x, y) \in \text{SUPP}(Test) \setminus \text{SUPP}(Source)$ , it is usually a small value.

Suppose (x, y) between the unseen domain for testing and source doamins have been fully covered by  $\Lambda_S$ , then  $\eta_{T/S} = 0$  and there exists a finite upper bound  $\rho = \sup_{i,j \in [N_d]} \beta_q(D_i || D_j), \forall q > 0$  for any convex combinatorial domains  $D_i$  and  $D_j$ . 

$$R_{D_U}[h] \le \frac{1}{2} d_{D_U}(h) + \rho \cdot \left[ e_{D_{\bar{U}}}(h) \right]^{1 - \frac{1}{q}},$$
(45)

where  $\rho$  denotes the maximum domain divergence of source domains we could minimize, and  $e_{D_{\pi}}(h)$  represents that the empirical risks of source domains to be minimized.

Proposition 5.2 (Frequency Perspective Risk Bound on Unseen Time series Domain).Let H be a hypothesis space built from a set of source time series domains  $D = \{D_i\}_{i=1}^{N_d}$ . Suppose q > 0is a constant, for any unseen time series domain  $D_U$  from the convex hull  $\Lambda_D$ , we have its closest element  $D_{\bar{U}}$  related to source domains in  $\Lambda_D$ , i.e.,  $D_{\bar{U}} = \arg \min_{\pi_1, \dots, \pi_{N_d}} \beta_q(D_{\bar{U}} || \sum_{i=1}^{N_d} \pi_i D_i)$ . Then the risk of  $D_U$  on any label function  $h \in \mathcal{H}$  is,

$$R_{D_U}[h] \le \frac{1}{2} d_{D_U}(h) + \rho \cdot \left[ e_{D_{\bar{U}}}(h) \right]^{1 - \frac{1}{q}}, \tag{46}$$

where  $\rho = 2^{\frac{q-1}{q} \sup_{i,j \in [N_d]} RD_q(D_i || D_j)}$ ,  $d_D(h)$  and  $e_D(h)$  are ideal and empirical risk of domain D,

969  
970 
$$RD_{q}(D_{i}||D_{j}) = \frac{1}{q-1} \log \int \left[P_{F}(D_{i})\right]^{q} \left[P_{F}(D_{j})\right]^{1-q} da = \frac{1}{q-1} \log \left[\frac{\mu_{j}^{2(q-1)}}{\mu_{i}^{2q}} \cdot \frac{\sqrt{\pi}^{2n+1}n!!}{2^{2n+1}\sqrt{\gamma}^{n+1}}\right]$$
(47)

where 
$$\mu_i = E(D_i) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \sigma_k, \mu_j = E(D_j) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \tau_k, \gamma = \sum_{k=1}^n \frac{q}{2\sigma_k^2} + \frac{1-q}{2\tau_k^2}, \{\sigma_k\}_{k=1}^n$$
  
and  $\{\tau_k\}_{k=1}^n$  denote frequency scale parameters in  $D_i, D_j, n$  represents the number of components.

*Proof.* Suppose the overall distribution supports condition is  $SUPP(D_i) = SUPP(D_i)$ , their corre-sponding overall probability distributions on the frequency domain  $P_F(x)$  also satisfy the condition, then it is feasible for us to use the generalized Rényi Divergence to estimate its whole probabilistic density function  $f(a) = \prod_{k=1}^{n} f(a_k)$ , It is a reasonable extension for the individual frequency do-main components mentioned in Raincoat(He et al., 2023) do not satisfy the KL-Divergence in the range of  $a_k \in (0, +\infty)$ . We define the probability density functions  $f_i(a), f_i(a)$  corresponding to the two distributions  $D_i$  and  $D_j$ , 

$$f_i(a) = \prod_{k=1}^n \frac{a}{\sigma_k^2} \cdot \exp(-\frac{a^2}{2\sigma_k^2}), f_j(a) = \prod_{k=1}^n \frac{a}{\tau_k^2} \cdot \exp(-\frac{a^2}{2\tau_k^2})$$
(48)

where  $\sigma$  and  $\tau$  denotes the scale parameters in  $D_i$ ,  $D_j$  respectively, while n is the number of frequency components. Then we can calculate the formulation of  $\frac{f_i(a)}{f_i(a)}$ ,

$$\frac{f_i(a)}{f_j(a)} = \prod_{k=1}^n \frac{\frac{a}{\sigma_k^2} \cdot \exp(-\frac{a^2}{2\sigma_k^2})}{\frac{a}{\tau_k^2} \cdot \exp(-\frac{a^2}{2\tau_k^2})} = \prod_{k=1}^n \frac{\tau_k^2}{\sigma_k^2} \exp\left(-a^2(\frac{1}{2\sigma_k^2} - \frac{1}{2\tau_k^2})\right)$$
(49)

Following these, we can continue to derive  $RD_q(D_i||D_i)$ ,

$$RD_q(D_i||D_j) = \frac{1}{q-1} \log \int \left[f_i(a)\right]^q \left[f_j(a)\right]^{1-q} da$$
(50)

$$= \frac{1}{q-1} \log \int \left[\frac{f_i(a)}{f_j(a)}\right]^q f_j(a) da$$
(51)

$$= \frac{1}{q-1} \log \int \left[ \prod_{k=1}^{n} \frac{\tau_k^2}{\sigma_k^2} \exp\left( -a^2 (\frac{1}{2\sigma_k^2} - \frac{1}{2\tau_k^2}) \right) \right]^q \cdot \prod_{k=1}^{n} \frac{a}{\tau_k^2} \cdot \exp(-\frac{a^2}{2\tau_k^2}) da$$
(52)

 $= \frac{1}{q-1} \log \int \prod_{k=1}^{n} \frac{\tau_k^{2q-2}}{\sigma_k^{2q}} a^n \cdot \exp\left(-qa^2 \sum_{k=1}^{n} (\frac{1}{2\sigma_k^2} - \frac{1}{2\tau_k^2}) - a^2 \sum_{k=1}^{n} \frac{1}{2\tau_k^2}\right)$ (53)

$$= \frac{1}{q-1} \log \int \prod_{k=1}^{n} \frac{\tau_k^{2q-2}}{\sigma_k^{2q}} a^n \cdot \exp\left(-a^2 \sum_{k=1}^{n} (\frac{q}{2\sigma_k^2} + \frac{1-q}{2\tau_k^2})\right)$$
(54)

Let  $A = \prod_{k=1}^{n} \frac{\tau_{k}^{2}}{\sigma_{k}^{2q}}, B = \sum_{k=1}^{n} (\frac{q}{2\sigma_{k}^{2}} + \frac{1-q}{2\tau_{k}^{2}}),$  $BD_{\alpha}(D_{i}||D_{i}) = \frac{1}{\log A} \int_{-\infty}^{+\infty} a^{n} \cdot exp(-Ba^{2})da$ 

$$RD_{q}(D_{i}||D_{j}) = \frac{1}{q-1} \log A \int_{0} a^{n} \cdot exp(-Ba^{2})da$$
(55)  
$$= \frac{1}{q-1} \log A \frac{\sqrt{\pi}n!!}{(2\sqrt{B})^{n+1}}$$
(56)

1013 Let 
$$\mu_i = E(D_i) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \sigma_k, \mu_j = E(D_j) = \prod_{k=1}^n \sqrt{\frac{\pi}{2}} \tau_k, \gamma = B,$$
  
1014  
1015  $RD_q(D_i || D_j) = \frac{1}{q-1} \log \left( \left[ \prod_{i=1}^n \frac{\tau_k^{2q-2}}{\sigma_i^{2q}} \right] \cdot \left[ \frac{\sqrt{\pi}n!!}{(2\sqrt{\gamma})^{n+1}} \right] \right)$ 

$$RD_{q}(D_{i}||D_{j}) = \frac{1}{q-1} \log \left( \left| \prod_{k=1}^{n} \frac{\tau_{k}^{2q-2}}{\sigma_{k}^{2q}} \right| \cdot \left[ \frac{\sqrt{\pi}n!!}{(2\sqrt{\gamma})^{n+1}} \right] \right)$$
(57)

$$= \frac{1}{q-1} \log \left( \frac{\mu_j^{2q-2}}{\mu_i^{2q}} \cdot (\frac{\pi}{2})^n \cdot \frac{\sqrt{\pi}n!!}{(2\sqrt{\gamma})^{n+1}} \right)$$
(58)

1020  
1021  
1022
$$= \frac{1}{q-1} \log \left[ \frac{\mu_j^{2(q-1)}}{\mu_i^{2q}} \cdot \frac{\sqrt{\pi}^{2n+1} n!!}{2^{2n+1} \sqrt{\gamma}^{n+1}} \right]$$
(59)

It is easy to determine  $\lim_{n \to +\infty} \frac{\sqrt{\pi}^{2n+1} n!!}{2^{2n+1} \sqrt{\gamma}^{n+1}} \to +\infty$ , At the same time, we plotted its monotonicity for  $n \in \mathbb{N}^+$ ,  $1 \le n \le 100$ ,  $\gamma \in \mathbb{N}^+$   $1 \le n \le 20$  in Figure 6 (a). it usually shows a monotonically increasing trend.

Let's consider  $R_{D_U}[h]$  where  $\rho = 2^{\frac{q-1}{q} \sup_{i,j \in [N_d]} RD_q(D_i \| D_j)}, RD_q(\cdot \| \cdot) \ge 0$ 

$$R_{D_U}[h] \le \frac{1}{2} d_{D_U}(h) + \rho \cdot \left[ e_{D_{\bar{U}}}(h) \right]^{1 - \frac{1}{q}}, \tag{60}$$

1029 1030

1039

1040 1041

1042

1028

1031 1032 (1) q > 1, both  $\rho$  and  $\left[e_{D_{\bar{U}}}(h)\right]^{1-\frac{1}{q}}$  monotonically increasing  $\uparrow$ , if we reduce  $\rho$  and  $e_{D_{\bar{U}}}$ , The upper 1033 bound of  $R_{D_U}[h] \downarrow$ . We could decrease n and the log part will decrease.

1034 1035 1036 (2) 0 < q < 1, both  $\rho$  and  $[e_{D_{\bar{U}}}(h)]^{1-\frac{1}{q}}$  monotonically decreasing  $\downarrow$ , When we minimize the 1036  $e_{D_{\bar{U}}}, [e_{D_{\bar{U}}}(h)]^{1-\frac{1}{q}} \uparrow$ , so we should decrease the  $\rho$  to let the upper bound of  $R_{D_U}[h] \downarrow$ , it requires 1037  $RD_q(\cdot \| \cdot)$  larger, in this case,  $\frac{1}{q-1} < 0$  so we also need to decrease the log part by decreasing n.  $\Box$ 1038

B DATASET

**B.1** DATASET INFORMATION

We list detail introduction of the datasets we used with FEDNet:

- Spurious Fourier (Gagnon-Audet et al., 2022) dataset is designed to study the impact of spurious correlations in one-dimensional signals under distribution shifts. It involves binary classification tasks based on the frequency characteristics of the signals. Each signal is constructed from Fourier spectra with one low-frequency peak and one high-frequency peak. The dataset comprises different domains, which are 10%, 80%, and 90%, representing the correlation between the low-frequency signal and the label. In contrast, the high-frequency signal maintains a consistent 75% correlation with the label across all domains.
- HHAR (Gagnon-Audet et al., 2022) dataset is used to study human activity recognition across different smart devices, such as smartphones and smartwatches. This dataset includes five source domains, each containing data gathered from a different device. The goal is to train models that can generalize to unseen devices, effectively ignoring spurious information from complex signals.
- UCIHAR (Anguita et al., 2012) dataset captures daily activities of 30 volunteers aged 19 to 48 using mobile phone sensors. It features a sampling frequency of 50 Hz and contains 1,318,272 time series samples, each with 9 initial features. The classification task involves identifying one of six activities: walking, sitting, lying down, standing, going upstairs, and going downstairs. The dataset is organized into 5 domains based on participants, with each domain comprising data from 6 volunteers.
- UniMiB-SHAR (Micucci et al., 2017) dataset consists of activity data gathered from three mobile phone sensors at a sampling frequency of 50 Hz, involving 30 participants aged 18 to 60. These participants performed 17 detailed actions, including 9 everyday activities and 8 types of falls. For evaluation, the dataset is divided into 4 domains. It includes a total of 1,569 time series samples, each with 453-dimensional features derived from the three sensors.
- Opportunity (Chavarriaga et al., 2013) dataset contains data from 4 volunteers performing 18 daily activities in a home environment, such as opening and closing the dishwasher, refrigerator, drawers, and etc. It uses various inertial sensors to enhance the generalization of OOD data. The dataset is sampled at 30 Hz, resulting in a total of 869,387 time series samples, each with 77 features. It is divided into 4 domains for evaluation purposes.
- DSADS (Barshan & Yüksek, 2014) dataset consists of 19 activities collected from 8 subjects wearing body-worn sensors on 5 different body parts. It captures a variety of daily and sports activities, providing comprehensive data for human activity recognition research.
- PAMAP (Reiss & Stricker, 2012) dataset includes data on 18 activities performed by 9 subjects, each wearing 3 sensors. This dataset focuses on physical activities and is designed to aid in the development of models for recognizing a wide range of movements.
- EMG (Lobov et al., 2018) consists of 6 types of gestures with 8 channels recorded from 36 participants sampled at 200 Hz.
- EEG (Goldberger et al., 2000) is a single channel EEG dataset collected from 20 subjects to classify 5 sleep stages.

Table 3:	Dataset	statistics
----------	---------	------------

1081								
1082	Dataset	Shape	Classes	Domains	samples	Subjects	Sensors	Frequency
1083	Spurious Fourier	(50, 1)	2	3	12,000	-	-	-
1084	HHAR	(500, 6)	6	5	13,674	9	5	25 Hz
1004	UCIHAR	(125, 45)	6	5	1,318,272	30	1	50 Hz
1085	UniMiB-SHAR	(151, 3)	18	4	11,771	30	1	50 Hz
1086	Opportunity	(30, 77)	17	4	869,387	4	72	30 Hz
1000	DSADS	(125, 45)	19	4	1,140,000	8	5	25 Hz
1087	PAMAP	(200, 27)	18	4	3,850,505	9	3	100 Hz
1088	EMG	(200, 8)	6	4	33,903,472	36	1	200 Hz
1089	EEG	(3000, 1)	5	4	-	20	1	-

**B.2** DATA PROCESSING 

We will provide preprocessing code for all datasets that need to be processed. The preprocessing methods will be consistent with those mentioned in other works and with open-source code. The detailed preprocessing procedures for some datasets are as follows: 

DSADS Cross-Position. The DSADS dataset consists of 5 sensors positioned at torso (T), right arm (RA), left arm (LA), right leg (RL), and left leg (LL). Each sensor records 9-dimensional variables (x, y, z accelerometers; x, y, z gyroscopes; x, y, z magnetometers) representing the position in space. Originally, the data had dimensions of 125x1x45. We split it into 5 domains based on the positions, resulting in final data dimensions of 125x1x9. This dataset is used to study more challenging domain generalization issues across different body parts. 

**PAMAP.** For the PAMAP dataset, we followed the processing method of Diversify, selecting all sample records and categories. The original data records were used with a fixed window size of 200 and a window overlap ratio of 50%.

Initial Domain splitting We list detail initial domains of the datasets we used in Table 4. 

Dataset	Domains	Infomation
Spurious Fourier	3	{d=10%, d=80%, d=90%}
HHAR	5	{Nexus 4, Galaxy S3, Galaxy S3 Mini, LG watch, Gear watch
UCIHAR	5	{0,1,2,3,4,5}
UniMiB-SHAR	4	{1,2,3,5}
Opportunity	4	{S1,S2,S3,S4}
DSADS	4	$\{(0,1), (2,3), (4,5), (6,7)\}$
PAMAP	4	$\{(0,1,2,11), (3,5,6,9), (7,8,10,13), (4,12)\}$
EMG	4	$\{(0-8), (9-17), (18-26), (27-35)\}$
EEG	4	$\{(0,1,2,3,4), (5,6,7,8,9), (10,11,12,13,14), (15,16,17,18,19)\}$

Table 4: Initial domain information.

#### SUPPLEMENTARY EXPERIMENTAL RESULTS С

## C.1 FULL ABLATION RESULTS

Table 5: Ablation Settings of FEDNet.

1125	Model	Time-	Determir	nistic	Time-Stochastic	Classification
1126	widder	$\mathcal{L}_{ELOB}$	$\mathcal{L}_{con1}$	$\mathcal{L}_{con2}$	$\mathcal{L}_{sto}$	$\mathcal{L}_{cls}$
1128	w/o $\mathcal{L}_{con}$		X X	×		
1129	w/o $\mathcal{L}_{sto}$		x	,	×	1
1130	FEDNet-cross FEDNet-contrast		√ x	X		
1131	i Ebriet contrast	· ·	~	•	•	•

We conduct ablation study to verify the impact of each technique in FEDNet on performance, and the details of our setup are provided in Table 5. As shown in Table 6, it is evident that  $\mathcal{L}_{con}$  and 1134  $\mathcal{L}_{det}$  significantly affect the model's performance. For all datasets, after removing the constraints 1135 of the  $x_{det}$  part, the performance generally shows a significant decrease, which indicates that there 1136 is Domain Shift in this part of the features, and removing  $x_{det}$  shows a very large decrease in the 1137 performance on all datasets. On the other hand,  $\mathcal{L}_{sto}$  has a smaller impact on performance and even 1138 results in negative optimization in some cases (e.g., when the target domain is S3 on the Opportunity dataset), this is due to the low signal-to-noise ratio in the information of the  $x_{sto}$  features. Addi-1139 tionally, we can see that FEDNet outperforms its three variants in most results, demonstrating the 1140 soundness of our design. 1141

1142 Abnormal results in Opportunity dataset. due to the original time series sequence window length 1143 L = 30 being much shorter than that of other time series containing limited information, after applying the Fourier transform, only L/2 frequency components remain, and further decomposition 1144 results in only 2-3 stable frequency components extracted, making prediction difficult in w/o  $\mathcal{L}_{con}$ . 1145

1149							
1150	Dataset	Target	$ $ w/o $\mathcal{L}_{con}$	w/o $\mathcal{L}_{det}$	w/o $\mathcal{L}_{sto}$	FEDNet-cross	FEDNet-contrast
1151	Spurious Fourier	0	50.12	50.78	74.38	74.22	75.31
1152		0	83.18	22.23	96.34	99.18	99.48
1150		1	83.98	22.44	96.88	98.98	98.98
1100	HHAR	2	55.45	22.27	87.64	90.97	91.59
1154		3	45.65	17.17	60.56	58.18	67.87
1155		4	42.37	15.27	45.92	38.11	00.01
1156		1	41 72	17.88	76.82	83 74	80.13
4457	UCIHAR	2	70.96	18.18	96.48	97.15	98.53
1157		3	51.10	18.92	89.58	93.21	94.64
1158		4	54.63	17.76	97.79	98.34	99.02
1159		1	23.69	10.93	53.42	57.55	58.85
1160	UniMiB-SHAR	2	19.21	24.52	44.25	63.63	70.15
1100	child of the	3	18.75	10.52	68.09	70.06	71.05
1161		5	31.87	18.79	39.93	44.63	44.29
1162		S1	12.21	63.23	84.62	84.80	85.02
1163	Opportunity	S2	15.26	62.38	81.20	79.75	81.45
1104	11 2	53 54	29.34	51.42	/8.//	//.89	79.21
1104			21.15	55.59	61.20	81.50	01.90
1165		0	11.14	5.26	40.89	37.54	40.54
1166	DSADS Cross Position	2	13.66	5.20	36.25	35.13	30.05
1167	DSADS Closs I ositioli	3	9.52	5.20	33.28	32.16	33.93
1107		4	13.53	5.26	35.08	32.78	33.04
1168			46 00	17.01	71 77	64 84	73.00
1169		1	41.28	16.48	85.39	82.13	87.10
1170	EMG	2	41.38	16.44	78.82	71.53	79.66
1171		3	45.12	16.42	78.08	69.87	79.85
1171		0	84.60	5.26	88.20	90.74	92.80
11/2	DGADG	1	76.84	5.26	83.11	82.63	84.86
1173	DSADS	2	86.31	5.26	87.89	8/.36	93.24
1174		3	02.07	5.20	80.85	/1.00	87.71
1175		0	65.24	44.56	66.93	66.53	67.48
11/0	ΡΔΜΔΡ	2	35.02	37.52	22.01	47.95	35.27
1176	1710171	3	66.45	68.03	68.08	64.73	69.80
1177		-	1				

Table 6: Ablation study of FEDNet.

## 1178 1179

1146 1147 1148

1180

1181

## C.2 TOP AVERAGE AMPLITUDE RATIO EFFECT

1182 1183

We further conduct key experiments on the proportion  $\alpha$  of the acquired high average amplitude. 1184 1185 The results in Figure 7 show that we only need to extract 5-20% of the original frequency as a time-deterministic component to reach state-of-the-art performance, which suggests that the use 1186 of the frequency-domain component as prior information is robust and conducive to the model's 1187 generalization ability.



Figure 7: Hyper-parameter study of  $\alpha$  on three datasets. The X-axis represents  $\alpha\%$  of top average amplitude frequency components, and the Y-axis represents accuracy in time series classification.

C.3 TIME-STOCHASTIC FEATURE STUDY

We conducted experimental analysis on the features of the Time-stochastic component across three different dimensions, demonstrating that these features can serve as auxiliary information to effectively enhance the model's generalization ability, with potential for further optimization. 

**Patch condition number (feature level).** We analyzed the original patch matrix after dividing the Time-stochastic feature into patches. We used the condition number  $cond(x_{patch}) = \sigma_{max}/\sigma_{min}$ where  $\sigma_{max}$  and  $\sigma_{min}$  represent the maximum and minimum singular values of  $x_{patch}$  to evaluate this part of the features, as shown in Table 7. We found that using non-overlapping patches effec-tively reduces the condition number, ensuring the numerical stability of the original matrix. It allows the module learning more effectively and further speeding up the training process. 

Table 7: Time-stochastic patch matrix condition number with different mask ratio  $\alpha$ .

D		$\alpha = 0$	0.2	$\alpha = 0$	0.4	$\alpha = 0.6$		
Datase		stride = P/2	stride = P	stride = P/2	stride = P	stride = P/2	stride = P	
Spurious Fo	urier	377.80	14.57↓	1211.58	<b>41.70</b> ↓	200.43	24.53↓	
HHAR		57.31	93.70	6702.17	<b>5661.14</b> ↓	4149784.84	3688642.06	
UCIHA	2	7600.79	<b>30.31</b> ↓	159401	37.88↓	147293444	<b>161.05</b> ↓	
SHAR		47.13	<b>11.21</b> ↓	9288.10	17.402↓	830192	38.3774↓	
EMG		27.26	84.076	2065.86	<b>152.440</b> ↓	1561521	<b>12334.38</b> ↓	
OPP		151.68	<b>2.0032</b> ↓	59.24	<b>1.8594</b> ↓	81.89	<b>1.78146</b> ↓	
DSADS		351.83	<b>26.7469</b>	2260.30	<b>15.1602</b> ↓	3226488.87	<b>37.6795</b>	
PAMA	)	49.98	71.068	4921.72	<b>245.448</b>	2968754.15	26859.82	

Time series missing value study (data level). We further analysed a more realistic scenario, i.e., we tried to add different missing rates of values to the training domains. We found it will affect the extraction of time-deterministic features with Freqency Filter. In this case, the supplementation of Time-stochastic brought significant gains. Table 8 shows that retaining time-stochastic features can improve model's generalization when the invariant features are not sufficiently extracted.

Table 8: time-stochastic gains with different time series missing rate.

Dataset	missing r	ate = 20%	missing r	ate = 40%	missing r	ate = 60%	missing r	ate = 80%
	w/o L <sub>sto</sub>	FEDNet						
UCIHAR	94.52	98.55	94.52	97.41	64.55	70.60	69.74	74.06
UniMiB-SHAR	52.60	52.86	50.52	51.56	38.54	40.36	30.98	38.80
EMG	67.43	67.49	38.62	40.73	16.55	27.88	16.55	23.59

**EEG Long signal case study (domain level).** We investigated the OOD generalization for the classification of EEG signals with an extremely long sequence (L=3000). As shown in Figure 8, the spectrum maps produced by ultra-long sequences are more likely to contain unknown frequency distributions. When we use domain 2 as the target domain, our model is unable to detect frequency

#### variations specific to the high-frequency part from the training domain, which may also belong to invariant features.

Figure 8: EEG training domain mean frequency spectrum.

This example shows us the shortcomings of FEDNet, i.e., the invariant features learnt from stable frequencies are not sufficient if we rely exclusively on the training domain to extract them, especially in the case of ultra-long sequences where many of the frequency components are lost. This is why we need to retain the time-stochastic module. From Table 9, there is a 3% difference between w/o  $L_{sto}$  and FEDNet, while general domain generalization methods are always the lowest. We have also tried to simplify the attention mechanism with the MLP and the deep separable convolution.

Table 9: EEG target domain 2 performance.

Method	Accuracy	F1-score	Precision	Recall
VREx	68.58	56.95	57.44	59.02
GroupDRO	69.40	55.83	58.31	56.53
ANDMask	69.44	57.47	59.40	57.97
FEDNetw/o Lsto	69.99	57.97	63.26	59.31
FEDNet+ self-attention	72.90	62.16	61.74	63.01
FEDNet+ DwConv	70.21	57.78	56.33	62.09
FEDNet+ MLP	71.85	62.94	61.80	65.19

## C.4 Empirical Domain divergence study with $\alpha$ isolated level

We theoretically analyzed the changes in the maximum distance of the dataset under four retention ratios:  $\alpha = [0.2, 0.4, 0.6, 0.8]$ . From Figure 10 and 9 We found that the maximum distance is generally smallest when  $\alpha$  is 0.2 or 0.4, and the overall trend shows that as the retention ratio decreases, the domain generalization distance also decreases. The only exception is the UniMiB-SHAR dataset, where the calculated denominator contains a very small value of  $\hat{\gamma} < 0.01$ , which causes the overall result to be significantly large. However, we conducted mask ratio experiments on UniMiB-SHAR and found that the optimal frequency is also concentrated at 0.2-0.4 From Table 13.



Figure 9: empirical max domain divergence with different  $\alpha$  when q = 2.

## C.5 INTRINSIC RELATIONSHIP STUDY OF ORTHOGONAL FREQUENCY DECOMPOSITION

1294 In addition to the orthogonal FFT, We conducted experiments using orthogonal wavelet mother func-1295 tions and non-orthogonal wavelet mother functions in Table 10 and found that orthogonal wavelet functions generally outperform the non-orthogonal case, which is consistent with our theoretical



Figure 10: empirical max domain divergence with different  $\alpha$  when q = 1/2.

insight Eq. (39). It is the orthogonality of the frequency components that allows us to treat the probability distributions of each frequency component as independent of each other. This independence allows us to model the effects of these distributions separately.

Table 10: Performance on different wavelet mother functions for frequency decomposition.

Datasat	Targat	db2	ortho	gonal	1	non-ortho	gonal
Dataset	Taiget	Accuracy	F1	Accuracy	F1	Accuracy	.5 F1
Spurious Fourier	d=10%	18.00	17.81	33.34	33.33	22.38	22.37
	0	97.49	97.20	97.30	97.02	96.12	95.75
	1	97.27	96.89	96.97	96.58	95.23	94.84
HHAR	2	88.99	88.08	89.30	88.59	89.20	88.42
	3	59.22	55.80	59.74	57.32	61.43	59.82
	4	43.92	53.52	40.88	50.70	41.89	50.59
	0	75.50	76.80	78.38	80.45	78.09	80.78
	1	69.86	67.61	72.84	70.10	71.52	67.75
UCIHAR	2	90.61	90.49	89.44	88.97	89.73	89.49
	3	80.12	79.27	77.60	75.58	75.39	72.64
	4	90.06	90.35	90.06	89.88	92.38	92.56
	1	62.50	50.78	64.58	52.57	62.50	50.14
UniMiD SHAD	2	59.86	47.92	60.72	55.72	60.54	51.65
UIIIMID-SHAK	3	75.65	50.35	74.67	45.85	73.02	43.85
	5	42.61	29.52	44.29	28.04	41.94	29.09
	0	82.86	56.23	83.82	58.06	83.53	55.42
Opportunity	1	81.41	51.73	81.53	54.27	81.93	55.75
Opportunity	2	74.14	36.99	76.32	39.95	74.52	36.25
	3	81.00	47.96	80.60	51.67	80.67	47.25

## C.6 TIME-DETERMINISTIC GAIN PHENOMENON FOR IRM

We observed that VREx performs the worst in this temporal OOD scenario. This is mainly because the principle behind IRM series methods is based on the assumption of domain feature invariance. Studies have shown that these methods typically fail when there is a significant marginal shift in the data itself.we selected representative IRM-related methods incorporated the Frequency Filter mod-ule. We use only the separated time-deterministic components to complete the OOD task. As shown in Table 12, we find that using time-deterministic features as a prior resulted in stable improvements across various IRM variants and it fits our theoretical insights in Appendix A. We verified that fil-tering Time-Deterministic features learnt with the help of IRM correlation yields smaller invariant domain distances than learning the original features directly by A-distance in Table 11. 

Table 11: $\mathcal{A}$ -distance on	IRM-based invariar	it features on DSADS.
--------------------------------------	--------------------	-----------------------

Dataset Target		IRM		IB-IRM		VI	REx	Ι	IB
Dataset	Taiget	full	$\alpha_{0.2/0.4}$	full	$\alpha_{0.2/0.4}$	full	$\alpha_{0.2/0.4}$	full	$\alpha_{0.2/0.4}$
	0	0.8338	0.8235↓	0.9329	0.8235↓	0.8524	0.8235↓	0.8833	<b>0.8235</b> ↓
	1	0.8482	0.8513	0.8679	<b>0.8390</b>	0.86377	<b>0.8307</b>	0.9195	0.8235
DSADS	2	1.0319	<b>0.8421</b> ↓	0.9391	0.9473	0.8235	0.8431	1.0061	0.8338
	3	0.8235	0.8534	0.9287	<b>0.9102</b> ↓	0.9752	0.8235↓	0.8482	0.8235
	avg.	$0.8841 \pm 0.09$	$0.8425 \pm 0.01$	$0.9171 \pm 0.03$	$0.8800 \pm 0.05$	$0.8787 \pm 0.06$	$0.8302 \pm 0.01$	$0.9142 \pm 0.06$	0.8260 + 0

1351														
1352	Dataset	Target	full o	IRM =20%	α=40%	full	IB-IRN $\alpha = 20\%$	$M_{\alpha=40\%}$	full	VREx $\alpha = 20\%$	$\alpha = 40\%$	full	$\frac{\text{IIB}}{\alpha = 20\%}$	$\alpha = 40\%$
1354	Spurious Fourier	d=10%	48.84 4	8.97↑	<b>49.91</b> ↑	49.66	<b>51.34</b> ↑	50.81↑	49.66	<b>50.69</b> ↑	47.38	50.34	<b>51.34</b> ↑	50.16
1356			90.14	88.39	<b>90.33</b> ↑	85.23	83.10	<b>85.64</b> ↑	89.60	94.40↑	<b>95.28</b> ↑	85.42	85.17	<b>86.02</b> ↑
1357 1358	HHAR	$\begin{vmatrix} 1\\2\\3 \end{vmatrix}$	89.41 9 88.06 9	0.29↑	89.83↑ <b>90.65</b> ↑	86.17 80.63	84.97 82.81↑	<b>87.34</b> ↑ 81.88↑ 63.02↑	91.00 84.58 65.55	93.70↑ 90.03↑ 67.34↑	<b>95.02</b> ↑ 87.23↑	87.37 85.10	87.49↑ 85.72↑	87.25 84.79
1359		4	52.03 <b>5</b>	31.02 3.38↑	53.38↑	42.57	62.88⊺ 43.24↑	03.92⊺ 43.58↑	65.55 54.05	53.04	69.90 56.42↑	48.99	64.10 49.66↑	04.02 53.04↑
1360 1361		0	93.95 <b>9</b> 67.55 <b>8</b>	5.97↑ 5 <b>.43</b> ↑	<b>97.69</b> ↑ 61.92	95.97 61.26	<b>99.71</b> ↑ 74.89↑	97.98↑ 72.41↑	89.34 66.23	96.83↑ <b>86.09</b> ↑	<b>98.85</b> ↑ 62.25	98.85 79.14	96.83 <b>87.75</b> ↑	95.68 82.12↑
1362	UCIHAR	$\begin{vmatrix} 2\\ 3 \end{vmatrix}$	98.83 <b>9</b> 95.27	<b>9.71</b> ↑ 94.64	98.83 <b>97.79</b> ↑	98.83 96.21	<b>99.71</b> ↑ 95.27	99.71↑ <b>97.79</b> ↑	97.65 83.28	<b>99.71</b> ↑ 85.43↑	98.53↑ <b>86.38</b> ↑	99.41 90.22	96.48 94.01↑	95.60 <b>95.27</b> ↑
1363 1364		4   S1	89.40 <b>9</b>	8.01↑ 5.99↑	96.03↑ 62.49↑	83.44	88.41↑ 57.76↑	88.74↑ 58.74↑	70.86	89.34↑ 55.12↑	85.70↑ 52.60	96.36 72.22	96.69↑ 72.23↑	89.07 52.17
1365	Opportunity	S2 S3	43.31 <b>4</b> 35.17 <b>3</b>	4.38† 9.06†	<b>44.43</b> ↑ 38.01↑	47.61 44.09	<b>47.71</b> ↑ 40.28	45.84 43.31	37.17 36.77	47.15↑ <b>40.86</b> ↑	<b>50.76</b> ↑ 40.41↑	71.47 71.47	53.15 33.44	<b>71.98</b> ↑ 37.29
1367		S4	46.07 5	2.08	<b>54.10</b> ↑	48.52	52.06↑	<b>54.01</b> ↑	46.41	53.48↑	<b>59.78</b>	44.70	51.02↑	44.37
1368 1369	EMG		84.56	70.60 83.24	71.65 <b>84.84</b> ↑	69.31 86.44	71.48↑ 87.76↑	70.13↑ 87.71↑	70.25	71.77↑ 84.12	7 <b>1.89</b> ↑ 84.23	64.85 77.18	68.19↑ 79.44↑	66.84↑ 76.79
1370		$\begin{vmatrix} 2\\ 3 \end{vmatrix}$	77.20 7	7.81¶ 7.67↑	78.95↑ 77.97↑	80.44 80.09	79.96 79.92	80.74 80.57	75.36	76.79	78.47↑ 77.20↑	75.84 62.08	71.89 71.35↑	/3.68 <b>74.42</b> ↑
1371 1372	DSADS	0	83.63 71.77 <b>7</b>	83.56 ′ <b>5.98</b> ↑	<b>85.21</b> ↑ 71.68	86.22 86.99	83.83 <b>88.09</b> ↑	83.52 86.38	83.27 75.56	83.41↑ 73.80	78.86 <b>75.83</b> ↑	79.04 80.70	90.04↑ 82.19↑	$\begin{array}{c} 88.46 \uparrow \\ 81.27 \uparrow \end{array}$
1373	DSADS	$\begin{vmatrix} 2\\ 3 \end{vmatrix}$	83.54 8 78.57 <b>8</b>	9.83↑ <b>3.36</b> ↑	<b>90.45</b> ↑ 78.59↑	86.98 80.29	84.43 <b>85.70</b> ↑	<b>90.32</b> ↑ 84.12↑	88.27 77.81	87.56 83.16↑	<b>90.86</b> ↑ 82.81↑	85.22 74.61	83.95 <b>75.48</b> ↑	<b>90.18</b> ↑ 70.13
1374 1375		0	63.44 <b>6</b> 50.31 <b>5</b>	3.69↑ 1.55↑	62.47 51.35↑	60.04 49.81	<b>63.00</b> ↑ 50.26↑	61.54↑ <b>51.26</b> ↑	62.88 54.88	63.85↑ 56.30↑	62.93↑ 54.52	63.21 57.39	63.66 <sup>↑</sup>	62.74 53.57
1376 1377	PAMAP	$\begin{vmatrix} 2\\ 3 \end{vmatrix}$	25.68 62.62 <b>6</b>	22.97 2.68↑	23.96 <b>63.71</b> ↑	26.01 62.97	26.03 63.82	25.21 60.22	22.68 62.10	<b>24.98</b> ↑ 62.15↑	23.80↑ <b>62.49</b> ↑	22.87 65.79	23.79 <sup>†</sup> 66.16 <sup>†</sup>	<b>27.11</b> ↑ <b>68.97</b> ↑
1378			60.16	53.91	58.07	59.11	55.73	58.85	55.99	53.39	<b>59.64</b> ↑	60.42	55.47	57.81
1379 1380	UniMiB-SHAR		69.41 ( 36.91	+3.40	52.00 67.11 38.26 <sup>+</sup>	66.45 42.28	40.54 54.61 40.60	57.12 59.87 42.62 <sup>+</sup>	63.16 41.28	43.11 60.20 38.93	55.09 67.11↑ 40.60	65.79 46 64	49.74 66.45↑ 44.63	<i>52.14</i> <b>71.05</b> ↑ <i>44.30</i>
1381		1 5	50.71	50.50	50.20	12.20	10.00	-12.02	1.1.70	50.75	10.00	10.04	1-1.05	1.50

## Table 12: Time-Deterministic Enhancement for IRM-based variants.

## C.7 VISUALIZATION STUDY

We provide some t-SNE and FFT masking spectrum visualizations as shown in Figure 11 and 12. We have chosen several methods for comparison: FreTS and Diversify, one ablation version FEDNet w/o  $L_{con}$ , and FEDNet. From the results, it can be seen that our proposed FEDNet has a more compact potential representation and the division between different labels is clear. In addition, Diversify's division between labels is clear, but the representation is scattered, potentially causing confusion from redundant information and random noise.



Figure 11: t-SNE visualizations on the UCIHAR dataset with target domain 0. The X-axis represents the first dimension and the Y-axis represents the second dimension. The colors denote class labels.



**Figure 12:** Figures (a) and (b) show the spectrograms before and after masking for domain 0 of the SHAR dataset. It can be observed that the retained frequencies are not exclusively low-frequency signals. Figures (c) and (d) display the polar plots before and after masking between domain 0 and domain 3 of the HHAR dataset. It can be seen that, through masking, the diversity between domains is more pronounced, while many redundant frequencies are eliminated.

Table 13: Performance on dataset UniMiB-SHAR with different  $\alpha$  ratio level.

Torgat		FI	T			D١	NТ	
Target	0.2	0.4	0.6	0.8	0.2	0.4	0.6	0.8
0	57.29	61.19	58.59	58.59	55.20	56.77	59.89	58.07
1	67.75	65.00	67.06	62.55	58.83	73.58	66.72	66.03
2	67.10	69.73	65.13	68.42	75.65	75.65	74.01	73.68
3	35.23	42.95	38.25	38.92	40.26	44.29	40.26	38.59

## 1427 C.8 MODEL COMPLEXITY

Given a time series data  $L \times C$ , where L denotes its input length and C represents the number of channels, FFT complexity is  $O(L \log L)$ . The time-deterministic complexity is  $O(2 \times N \times L \times C \times C_{out} \times K)$  comes from the encoder and decoder, where N is the number of hidden layers, K is the kernel size of the 1D-Conv and  $C_{out}$  is the output channels. The time-stochastic part uses patchify to reduce tokens in attention layer, with complexity  $O(M^2)$  where  $M = L/S \ll L$ , S is patch stride.

C.9 FULL METRIC RESULTS

In order to more fully illustrate the effectiveness of our method, we provide results for other metrics in table 14, and it can be seen that our method achieves the best performance on most of the datasets.

# 1442 D IMPLEMENTATION DETAILS

D.1 EXPERIMENTAL ENVIRONMENT

We implement FEDNet and the baselines based on WOODS(Gagnon-Audet et al., 2022) Benchmark and Time-Series-Library on a server equipped with an Intel(R) Xeon(R) Gold 5117 CPU and a Tesla
V100 (32 GB) GPU, with 256 GB of memory. The server runs on Ubuntu 18.04 with CUDA 12.4, and the codes are implemented in PyTorch 2.3.0+cu121, Python 3.9.12. To reduce randomness, we conducted each experiment 3 times and reported the best results.

D.2 COMPARISON METHOD INTRODUCTION

We list detail introduction of the methods we compared with FEDNet:

IRM (Arjovsky et al., 2019) proposes a framework named Invariant Risk Minimization that aims to find representations where the optimal classifier remains invariant across different environments, thereby improving generalization to unseen domains.

Dataset	Target	VREx	GroupDRO	ANDMasl	k FreTS	PatchTST	Г GILE	AdaRNN	Diversify	FEDNet <sub>f</sub>	FEDNet <sub>w</sub>
Spurious Fourier	d=10%	32.52	32.73	32.93	11.15	50.32	11.03	<u>33.39</u>	15.37	74.56	33.33
HHAR	$\begin{vmatrix} 0\\1\\2\\3\\4 \end{vmatrix}$	88.57 90.31 83.61 60.32 56.16	88.24 86.66 82.80 <u>64.07</u> <u>59.49</u>	90.73 91.87 85.83 61.13 49.74	92.25 93.00 88.67 58.42 54.47	78.25 78.60 74.68 36.53 45.60	93.75 96.57 87.33 62.67 49.70	90.73 68.58 67.69 46.30 41.63	91.82 91.99 86.87 51.67 49.29	96.63 97.96 89.75 61.53 59.59	<b>97.51</b> <u>96.89</u> <b>92.25</b> <b>65.94</b> 53.52
UCIHAR	$\left \begin{array}{c}0\\1\\2\\3\\4\end{array}\right $	89.23 63.33 97.51 83.18 68.91	91.70 51.40 96.81 81.81 60.30	<b>98.73</b> 63.00 <b>97.85</b> 87.78 89.75	84.01 52.58 92.10 76.67 88.89	67.52 74.50 66.59 80.05 77.65	86.44 75.88 97.31 87.59 92.89	85.29 56.36 75.05 74.25 65.23	79.67 66.78 92.39 85.61 92.18	95.62 77.49 97.80 88.80 98.35	80.45 70.10 90.49 79.27 90.35
UniMiB-SHAR	$\begin{vmatrix} 1\\2\\3\\5 \end{vmatrix}$	43.37 51.37 31.59 23.70	45.99 53.55 35.52 24.74	45.89 52.55 35.69 23.56	27.96 32.41 41.62 23.10	30.85 54.49 51.71 24.85	33.80 47.09 48.20 <u>31.82</u>	29.10 17.75 15.73 <b>34.02</b>	36.21 28.66 40.50 20.10	$   \begin{array}{r}     41.05 \\     \underline{55.97} \\     \underline{49.81} \\     \underline{22.39}   \end{array} $	<b>52.57</b> <b>61.87</b> <b>50.35</b> 28.04
Opportunity	S1   S2   S3   S4	39.98 32.06 26.53 29.30	40.16 33.56 32.15 32.67	28.67 30.52 21.79 29.20	55.41   54.47   41.20   46.75	27.18 24.75 23.51 24.66	60.76 55.69 45.64 49.51	47.48 40.71 31.58 33.98	53.75 46.72 34.39 46.03	63.22 <u>55.66</u> <u>42.35</u> <b>51.79</b>	58.06 54.27 39.95 <u>51.67</u>
EMG	$ \begin{array}{c c} 0\\ 1\\ 2\\ 3 \end{array} $	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	70.55 82.71 <u>78.49</u> 78.85	71.25 82.10 76.60 <b>79.87</b>	71.57 80.16 74.62 78.26	32.07 35.66 14.35 36.14	62.97 67.97 65.88 70.02	52.58 56.39 56.04 52.76	65.18 78.61 70.63 76.50	72.71 87.06 79.46 77.69	64.03 59.57 77.07 <u>79.80</u>
DSADS	$\begin{vmatrix} 0\\1\\2\\3 \end{vmatrix}$	76.83 77.91 85.39 73.33	82.88 76.51 82.75 74.93	79.73 74.48 81.08 77.58	78.28 69.76 83.35 72.36	81.82 73.32 82.50 78.28	89.41 78.26 85.37 77.38	82.58 79.74 83.07 69.54	76.58 77.29 85.12 70.04	92.65 82.74 93.21 87.64	<u>92.57</u> <b>83.73</b> <u>90.27</u> <u>80.52</u>
PAMAP	$\begin{vmatrix} 0\\1\\2\\3 \end{vmatrix}$	45.22 46.14 22.70 50.60	48.18 41.24 26.21 52.77	47.49 43.83 27.71 <u>57.60</u>	48.79 51.49 35.50 <b>68.23</b>	48.66 55.46 <b>49.65</b> 51.54	<b>55.52</b> 41.45 26.52 55.06	53.61 43.92 23.13 47.57	53.78 41.31 25.40 48.67	<u>55.43</u> <b>59.90</b> 35.34 56.68	53.26 49.69 <u>37.15</u> 55.34

Table 14: Macro-F1 on cross-person generalization.

1458

1459

1490 1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

- IIB (Li et al., 2022) proposes a method to enhance domain generalization by extracting features with information bottleneck that are invariant across different domains, improving model robustness and performance.
- VREx (Krueger et al., 2021) introduces a method to enhance model robustness by penalizing the variance of risks across different environments, encouraging the model to perform consistently across diverse settings.
- IB-IRM (Ahuja et al., 2021) presents a method that combines the principles of empirical risk minimization with the information bottleneck approach to encourage models to focus on relevant features, enhancing robustness and generalization across different environments.
- GroupDRO (Sagawa\* et al., 2020) is a method that seeks a global distribution with the worst performance within a range of the raw distribution for better generalization. Ours study the internal distribution shift instead of seeking a global distribution close to the original one.
- ANDMask (Parascandolo et al., 2021) is another gradient-based optimization method that belongs to special learning strategies. Ours focuses on representation learning.
- GILE (Qian et al., 2021) is a disentanglement method designed for cross-person human activity recognition. It is based on VAEs and requires domain labels.
- AdaRNN (Du et al., 2021) is a method with a two-stage that is non-differential and it is tailored for RNN. A specific algorithm is designed for splitting. Ours is universal and is differential with better performance.
- Diversify (Lu et al., 2022) is a time series OOD generalization method for dynamic distribution representation learning. It constructs a set of latent domain labels to better adapt to downstream tasks by employing a min-max adversarial approach to divide the original time series data distribution.

Dataset	Target  VREx	GroupDRO	O ANDMas	k FreTS I	PatchTS	Γ GILE Δ	AdaRNN	Diversify	$FEDNet_f$	FEDNet
Spurious Fourier	d=10% 24.33	24.33	24.55	11.17	<u>50.32</u>	11.03	25.06	15.36	74.56	33.34
HHAR	$ \begin{vmatrix} 0 & 89.14 \\ 1 & 90.21 \\ 2 & 84.09 \\ 3 & 65.41 \\ 4 & 54.24 \end{vmatrix} $	88.57 86.79 83.29 65.45 56.35	90.84 91.76 86.17 61.90 48.61	92.41 92.99 89.13 63.63 57.18	79.14 78.44 75.07 36.48 46.99	94.39 96.45 88.03 <u>65.93</u> 51.45	90.84 77.17 68.03 46.46 <b>60.84</b>	91.14 91.55 84.79 54.32 51.49	96.65 97.98 89.43 62.82 59.85	<b>97.73</b> 97.32 <b>92.38</b> <b>69.86</b> 53.28
UCIHAR	0 91.19 1 79.53 2 97.86 3 88.69 4 82.45	93.04 74.97 97.24 89.01 62.54	<b>98.53</b> 61.07 <u>97.82</u> <u>91.13</u> <u>91.01</u>	87.24 69.66 92.47 83.49 90.26	78.13 75.30 71.86 81.09 80.89	89.70 83.72 97.31 90.44 94.07	89.29 54.34 80.91 81.06 68.94	80.78 67.39 93.14 87.85 93.27	95.36 86.99 97.80 92.64 98.36	80.89 79.52 90.79 81.24 91.53
UniMiB-SHAR	$ \begin{vmatrix} 1 \\ 2 \\ 3 \\ 54.57 \\ 35.65 \\ 5 \\ 24.03 \end{vmatrix} $	$\begin{array}{r} 48.37 \\ 57.70 \\ 42.31 \\ 24.70 \end{array}$	49.15 57.22 37.17 24.21	28.49 35.72 40.15 24.84	38.93 54.74 49.69 23.53	33.11 46.52 55.08 29.04	35.86 18.06 13.63 17.85	39.31 32.61 43.10 <b>36.41</b>	43.21 57.58 57.77 25.63	<b>52.90</b> <b>63.03</b> <u>56.00</u> 23.83
Opportunity	S1         34.96           S2         28.09           S3         24.86           S4         25.20	34.87 28.76 30.01 27.36	39.85 45.41 26.85 39.40	63.66 61.58 52.73 45.92	23.50 24.75 20.60 25.36	$\frac{68.60}{62.78} \\ \frac{55.80}{59.30} \\$	58.21 55.83 43.54 43.99	59.25 60.05 47.11 53.81	70.79 64.57 56.41 59.96	68.09 <b>64.84</b> 51.36 58.61
EMG	0 72.43 1 84.37 2 73.51 3 78.27	71.12 84.61 78.66 79.48	71.92 84.21 77.30 <b>80.05</b>	72.52 81.49 74.99 78.95	31.48 35.34 10.81 37.74	64.20 69.04 66.04 70.56	53.48 56.88 56.80 54.84	66.91 80.67 71.28 76.79	<b>73.19</b> <b>87.31</b> <b>80.32</b> 79.39	67.19 60.34 78.00 78.43
DSADS	$ \begin{array}{c ccccc} 0 & 81.79 \\ 1 & 87.21 \\ 2 & 90.89 \\ 3 & 79.48 \end{array} $	88.22 82.31 85.59 81.53	86.59 82.16 84.49 81.54	82.35 77.73 88.50 80.42	84.00 76.07 86.29 81.81	91.23 84.72 91.47 <u>85.30</u>	85.26 83.58 89.24 75.67	80.33 83.51 89.83 77.72	93.65 85.51 93.75 88.67	<b>93.93</b> <b>88.57</b> <u>92.58</u> 83.35
PAMAP	0 44.60 1 56.77 2 25.36 3 50.76	49.92 49.20 29.64 54.28	47.62 50.45 30.35 <u>62.41</u>	53.45 52.49 33.46 <b>73.39</b>	53.78 57.26 <b>49.55</b> 55.34	56.25 49.01 30.11 55.41	<b>59.30</b> 44.12 27.57 48.42	57.37 42.01 26.73 50.61	58.89 64.78 42.35 59.73	57.37 58.01 39.23 55.53

Table 15: Precision on cross-person generalization

1513

1538 1539

1541

#### D.3 EVALUATION 1540

For UCIHAR, UniMiB-SHAR, and Opportunity, we directly utilized the open-source original 1542 datasets processed by GILE. To faithfully replicate the performance of GILE on these datasets, 1543 our experimental evaluation follows the same methodology as GILE. Specifically, for each domain, 1544 we treat it as the testing set while the remaining domains serve as the training set. We then select 1545 the model with the highest accuracy on the testing set. 1546

For Spurious Fourier and HHAR, both are derived from the Benchmark dataset provided by 1547 WOODS (Gagnon-Audet et al., 2022). We adopt WOODS' default methods and domain parti-1548 tioning as our baseline. This involves splitting 20% of the data in each domain for evaluation, using 1549 the evaluation data from the training domain as the validation set, and the evaluation data from the 1550 testing domain as the testing set. We calculate the average accuracy of each checkpoint on all vali-1551 dation sets across domains and save the results corresponding to the model with the highest average 1552 accuracy. The evaluation metric primarily relies on the Train-domain validation from WOODS, as 1553 it aligns with the evaluation methodology used in other datasets. 1554

For DSADS, PAMAP, and we improved the architecture provided by WOODS and extended it to 1555 integrate these datasets for evaluation. The dataset partitioning method for evaluation follows the 1556 same approach as employed in Diversify, where only 20% of the training data is set aside as the 1557 validation set, and the remaining data from each domain is used as the testing set. We record the 1558 results on the testing set corresponding to the model that performed best on the validation set. 1559

1560

#### 1561 D.4 HYPERPARAMETER SETTING

1562

For hyper-parameter settings, if baselines provide hyper-parameters for the used datasets, we keep 1563 their default settings. Otherwise, we adjust the hyper-parameters to ensure a fair comparison as 1564 much as possible. For our proposed FEDNet, we leverage a 4-layer CNN with max pooling for the 1565 time-deterministic block and a 2-layer transformer encoder for the time-stochastic block, attention

					-	•				
Dataset	Target  VREx (	GroupDRC	) ANDMas	k FreTS I	PatchTS	Γ GILE /	AdaRNN	Diversify	FEDNet <sub>f</sub>	FEDNet
Spurious Fourier	d=10% 50.00	50.00	50.00	11.15	50.32	11.03	50.00	15.38	74.56	33.33
HHAR	$ \begin{vmatrix} 0 & 88.38 \\ 1 & 90.79 \\ 2 & 83.72 \\ 3 & 64.30 \\ 4 & 64.50 \end{vmatrix} $	88.39 86.96 83.02 63.99 <u>67.58</u>	90.89 92.08 86.01 61.40 56.46	92.22 93.33 88.46 59.73 61.40	78.48 79.21 74.87 38.71 50.52	93.53 96.75 87.74 62.29 55.90	90.89 71.20 70.48 48.36 52.32	91.20 91.54 83.41 53.42 59.33	96.64 97.94 90.78 61.86 69.21	<b>97.34</b> 96.57 <b>92.17</b> <b>66.79</b> 58.95
UCIHAR	$ \begin{vmatrix} 0 & 88.65 \\ 1 & 66.59 \\ 2 & 97.35 \\ 3 & 82.92 \\ 4 & 72.84 \end{vmatrix} $	92.02 57.98 96.79 84.39 67.48	<b>98.96</b> 67.74 <b>97.94</b> <u>88.80</u> 90.15	85.37 56.82 91.92 79.22 89.53	76.12 72.67 70.84 78.89 80.18	88.14 75.88 97.30 87.59 92.89	85.53 61.43 75.49 73.09 69.62	80.01 70.99 92.26 85.73 92.05	96.66 78.97 97.81 91.72 98.39	80.35 72.15 90.35 78.95 90.61
UniMiB-SHAR	$ \begin{vmatrix} 1 \\ 2 \\ 3 \\ 55.12 \\ 36.90 \\ 5 \\ 30.95 \end{vmatrix} $	$\frac{47.81}{56.41}\\39.58\\33.12$	49.19 56.63 40.78 31.57	28.49 36.32 46.51 29.50	$\begin{array}{r} 33.54 \\ \underline{60.54} \\ 54.62 \\ 31.54 \end{array}$	39.79 52.94 52.14 38.41	31.75 25.34 21.75 19.97	37.96 40.04 44.40 40.91	45.61 60.39 49.33 27.31	<b>54.84</b> <b>64.89</b> <b>50.68</b> 37.19
Opportunity	S1 <b>59.88</b> S2 <u>53.15</u> S3 <b>43.82</b> S4         48.42	56.63 48.13 <u>42.81</u> 51.01	25.86 25.95 19.74 28.10	50.77 50.30 38.82 45.92	39.82 27.77 33.52 35.68	55.96 <b>54.19</b> 42.28 49.83	43.19 35.38 29.24 32.75	50.60 41.20 34.73 43.42	<u>58.95</u> 53.07 38.31 <u>51.10</u>	52.03 49.94 36.80 <b>51.55</b>
EMG	$ \begin{array}{c cccc} 0 & 72.02 \\ 1 & 83.28 \\ 2 & 73.21 \\ 3 & 77.16 \end{array} $	70.75 83.19 <u>78.55</u> 78.64	71.69 82.61 76.62 79.86	71.87 80.25 74.59 78.02	34.37 36.94 22.55 36.46	63.21 68.26 65.77 70.12	54.32 57.59 57.39 53.82	66.23 79.13 71.11 76.43	<b>73.11</b> <b>87.17</b> <b>79.39</b> 77.54	64.73 59.98 77.14 <b>79.97</b>
DSADS	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	84.69 76.58 84.69 76.62	82.50 73.42 83.03 78.46	80.26 70.13 84.29 73.46	82.24 74.07 82.67 78.85	89.64 78.20 86.75 79.56	83.11 79.78 83.46 70.35	77.19 77.28 85.22 71.80	92.80 84.86 93.24 87.71	$\begin{array}{r} 92.41 \\ 83.64 \\ 90.65 \\ 80.52 \end{array}$
PAMAP	$ \begin{vmatrix} 0 & 49.10 \\ 1 & 45.03 \\ 2 & 27.03 \\ 3 & 51.47 \end{vmatrix} $	52.46 41.99 30.12 52.72	52.46 43.18 30.72 56.43	51.64 52.63 41.72 <b>70.00</b>	48.71 54.89 <b>55.83</b> 50.93	58.56 42.09 30.26 57.91	<b>58.59</b> 51.22 25.94 47.65	57.79 44.97 27.10 49.58	57.85 <b>61.34</b> 36.56 <u>58.11</u>	56.66 49.73 40.06 57.79

Table 16	: Recall	on	cross-person	generalization.
				<b>U</b>

1566

layer with 8 multi-head, patch length set to 16. We set the temperature for contrastive learning to
 0.07 or 0.2.

For UCIHAR, UniMiB-SHAR, Opportunity, we follow the same settings with GILE The hidden dimension is set to 50 for UCIHAR and UniMiB-SHAR datasets, and 128 for Opportunity datasets as they are more complex.

We used the *Adaptive Moment Estimation* (Adam) optimizer for all our training processes, with learning rates primarily adjusted 1e-2 ~1e-5. The weight decay parameter is typically set to {0, 1e-5, 5e-4}. For the UCIHAR, UniMiB-SHAR, and Opportunity datasets, we employed the WeightedRandomSampler consistent with GILE to balance label distribution. For other datasets without specific requirements, we utilized the RandomSample method provided by WOODS (Gagnon-Audet et al., 2022) for class balancing strategies.

1604 1605 1606

Table 17: Training hyperparameter settings for specific objective.

Object	Hyperparameter	Value	Object	Hyperparameter	Value
VREx	penalty weight annealing iterations	1e4 [500,1000,2000,4000]	Diversify	latent_domains alpha	[5,10] [0.1,1.0,10]
GroupDRO	$\eta$	1e-2	ANDMask	au	1.0

1611 1612

## <sup>1613</sup> E FURTHER DISCUSSION

1614

Finer-grained analysis with temporal stochastic components. Our method effectively separates the influence of two types of shifts and models the relationship between time-stable modules and domain shifts well. We have verified that this approach can effectively mitigate the generalization performance of IRM methods under data conditions where the marginal probabilities within domains dynamically change over time. However, finer-grained research on the impact of temporal shift components still requires further investigation. From the current experimental results, better

1621			51	
1622		<b>D</b>		
1623	Objective	Dataset	Model	Model-parameter
1624	VDEv			
1625	VKEX	Spurious		hidden_dept=3,
1626	GroupDRO	Fourier	LSTM	recurrent_lavers=2.
1627	ANDMost			state_size=32
1628	ANDMask			
1629	VRFv			
1630	VICEX			n_filters_time=32,
1631	GroupDRO	HHAR	Deep4Net	n_filters_spat=32,
1632	ANDMost			$n_{\rm h}$ nters = [64,128,256]
1633	ANDMASK			
1634	VRFx	DSADS		
1635	VILLA	PAMAP		kernel_size=(1,6)
1636	GroupDRO	UCIHAR	ActNetwork	MaxPool2d_kernel_size=(1,2)
1637	<b>ANDM</b> ask	Opportunity		MaxPool2d_stride=2
1638	AITDMASK			
1639			FreTS	embed size=128, hidden size=256
1640	ERM			
1641			PatchTST	$patch_len=16$ , $d_model = 128$ , $n_heads = 4$
1642	6H 5			
1643	GILE	ALL	GILE	kernel_size=9, d_AE=50
1644	AdaRNN		AdaRNN	n_hiddens=[64,64], trans_loss='mmd'
1645				
1646	Diversify		ActNetwrok	kernel_size=6, alpha1=1.0, alpha=1.0
1647	EEDN-4		EEDN-4	kernel_size=9, hidden_size=50,
1648	rednet		FEDNet	patch_len=16, d_model=512, n_heads=8

Table 18: Model hyper-parameter settings.

modeling of dynamic changes is beneficial for improving model robustness. Currently, there are
 related works on large language models attempting to map the changes in time series to state tokens
 that describe temporal trends. In the future, we can use the results of temporal shift components as
 guideline to model trends between patches, capturing the patterns of different domains change over
 time.

multi-domain datasets training phenomenon. We observed a phenomenon in the training of multi-domain datasets, which mainly arises from differences in the training processes of diversify and GILE compared to WOODS. For the training processes of diversify and GILE on UCIHAR, UniMiB-SHAR, and Opportunity datasets, 20% of the data is randomly divided from the entire training data. This can lead to an imbalance in the number of different training domains. In contrast, WOODS uses a method where 20% is divided from each domain as a validation set, and the multi-domain data is simultaneously loaded for training. This training strategy is more balanced compared to the previous method, but it also reduces the model's performance because, after balancing each domain, the originally smaller number of samples becomes even fewer. We hope to discuss which of these two training strategies is more reasonable, or if there is a better evaluation, as this will be helpful for our subsequent research.