

---

# Invariant Spatiotemporal Representation Learning for Cross-patient Seizure Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Automatic seizure type classification from electroencephalogram (EEG) data can  
2 help clinicians to better diagnose epilepsy. Although many previous studies have  
3 focused on the classification problem of seizure EEG data, most of these methods  
4 require that there is no distribution shift between training data and test data, which  
5 greatly limits the applicability in real-world scenarios. In this paper, we propose an  
6 invariant spatiotemporal representation learning method for cross-patient seizure  
7 classification. Specifically, we first split the spatiotemporal EEG data into different  
8 environments based on heterogeneous risk minimization to reflect the spurious  
9 correlations. We then learn invariant spatiotemporal representations and train  
10 the seizure classification model based on the learned representations to achieve  
11 accurate seizure type classification across various environments. The experiments  
12 are conducted on the largest public EEG dataset, the Temple University Hospital  
13 Seizure Corpus (TUSZ) dataset, and the experimental results demonstrate the  
14 effectiveness of our method.

## 15 1 Introduction

16 Epilepsy is a pervasive neurological disease that affects individuals all over the world, with consid-  
17 erable cognitive, psychological, and social ramifications [4]. The mainstream approach to epilepsy  
18 diagnosis relies on EEG data to classify seizures [8, 9]. However, traditional methods based on  
19 human labor are not only costly, but also susceptible to human uncertainty, as these methods require  
20 clinicians to meticulously review extensive EEG recordings [17]. As a result, using machine learning  
21 techniques to automatically classify seizure type attract increasingly attentions.

22 Early machine learning methods for accurately classifying EEG data included Support Vector Ma-  
23 chines (SVM), k-Nearest Neighbors (k-NN), and Bayesian methods [19, 32]. With the advancement  
24 of deep learning, Convolutional Neural Networks (CNNs) [34] and Recurrent Neural Networks  
25 (RNNs) [33] have been introduced. CNN-based methods typically aim at learning spatiotemporal  
26 feature representations of EEG signals through convolutional operations [6], exemplified by EEG-  
27 DBNet [24] and ACPA-ResNet [41]. RNNs, including CNN-BiRNN and CNN-Bi-LSTM [15, 25],  
28 capture temporal dependencies and dynamics. To address non-Euclidean geometric properties over-  
29 looked by CNNs and RNNs, Graph Neural Networks (GNNs) have been proposed to model the  
30 spatial relationships between EEG electrodes using a graph representation [10, 12, 18]. Methods such  
31 as REST [1], DCRNN [35], NeuroGNN [11] integrate GNNs with recurrent structures to enhance  
32 classification by capturing spatiotemporal dependencies and dynamic interactions.

33 However, these aforementioned methods are predominantly patient-specific and rely on a consistent  
34 distribution between training and test sets, which limits their ability to address cross-patient problem  
35 [40]. This kind of problems can be partially attributed to the spatial-temporal evolution of EEG  
36 data, which is common in real-world scenario where data from different patients exhibit significant

37 variability [22, 43]. Thus, for a group of new patients, it is very likely that this shift will impact the  
 38 performance of models, leading to less precision and reduced generalizability. These challenges  
 39 underscore the crucial and urgent need to develop robust cross-patient methods.

40 Previous methods addressing the cross-patient problem can be broadly categorized into three types.  
 41 The first type involves unsupervised representation learning, particularly domain generalization, to  
 42 initialize more robust representations for downstream tasks [38, 39, 42]. The second type revolves  
 43 around supervised models to enhance generalization by employing techniques such as causal learning  
 44 and invariant risk minimization. These approaches emphasize end-to-end learning strategies, which  
 45 have been shown to improve robustness to distributional shifts [26–28]. The third type involves  
 46 optimization-based approaches, including distributionally robust optimization (DRO), which focuses  
 47 on minimizing the worst-case performance under potential shifts in the data [21, 30]. However, most  
 48 of these methods ignore the spatiotemporal information, which leads to sub-optimal performance.

49 In this paper, we proposed a novel spatiotemporal invariant risk minimization loss to solve this  
 50 problem. Specifically, we first use the invariant mask function to separate the raw EEG feature into  
 51 the invariant representation and variant representation, and use the self-supervised learning (SSL) to  
 52 guarantee the preserved invariant information is able to predict the invariant feature at next timestamp.  
 53 In addition, we use the label information to generate the supervised signal to ensure the preserved  
 54 invariant information can predict the seizure type. Finally, we use the variance of the gradient toward  
 55 the mask function to control the time-varying variation of our methods in different patient groups.

56 We highlight our contributions as follows:

- 57 • We use the mask function to capture the invariant spatiotemporal information in the raw  
 58 EEG data and use such information for self-supervised learning.
- 59 • To further control the variation of the loss of the classification model, we use the variance of  
 60 the gradient as the penalty to achieve invariant learning.
- 61 • The experiments on the largest public dataset verify the effectiveness of our method.

## 62 2 Preliminary

### 63 2.1 Problem Setup

64 The primary objective of the seizure classification task is to predict the seizure type from a given  
 65 EEG signal clip. These clips were sliced from seizure EEGs using non-overlapping sliding windows  
 66 with fixed temporal size  $T$ . For each sample, we denote  $X \in \mathbb{R}^{T \times N \times M}$  as the EEG clip feature  
 67 after preprocessing, where  $T$  is the temporal length of the EEG clip,  $N$  is the number of EEG  
 68 channels/electrodes, and  $M$  is the number of features obtained through Fast Fourier Transform (FFT).  
 69 Meanwhile, we denote  $y$  as the seizure class label. For the independent identical distributed scenario,  
 70 different clips from the same patient may appear in both the training and test sets. However, in real  
 71 healthcare scenarios, patients in the test sets (a group of new patients) may completely unseen in the  
 72 training set, leading to the cross-patient problem [44], which can be further formulated as follows:  
 73 The patient set  $P$  is divided into two disjoint subsets,  $P_T$  and  $P_D$ , such that  $P_T \cup P_D = P$  and  
 74  $P_T \cap P_D = \emptyset$ . Here,  $P_D$  is used for training, and  $P_T$  is used for testing.

### 75 2.2 Previous Graph-Based methods for EEGs

76 **Graph Representing.** Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathbf{W}\}$  denote the graph structure, where  $\mathcal{V}$  is the set of nodes,  
 77  $\mathcal{E}$  refers to the set of edge, and  $\mathbf{W}$  is the adjacency matrix of the graph. In consideration of the  
 78 distribution of nodes and the physiological properties of the brain, two distinct approaches to graph  
 79 construction on EEG clips are evident. One undirected distance graph-based approach is to utilize the  
 80 Euclidean distance between different nodes on standard 10-20 EEG electrode placement as the basis,  
 81 followed by the threshold Gaussian kernel to determine the weights between  $v_i$  and  $v_j$  ( $v_i, v_j \in \mathcal{V}$ ):

$$W_{ij} = \begin{cases} \exp\left(-\frac{\text{dist}(v_i, v_j)^2}{\sigma^2}\right) & \text{if } \text{dist}(v_i, v_j) \leq \zeta \\ 0 & \text{otherwise,} \end{cases}$$

82 where  $\text{dist}(v_i, v_j)$  represents the Euclidean distance between two nodes  $v_i$  and  $v_j$ ,  $\sigma$  is the standard  
 83 deviation of the distances, while  $\zeta$  is the threshold for sparsity.

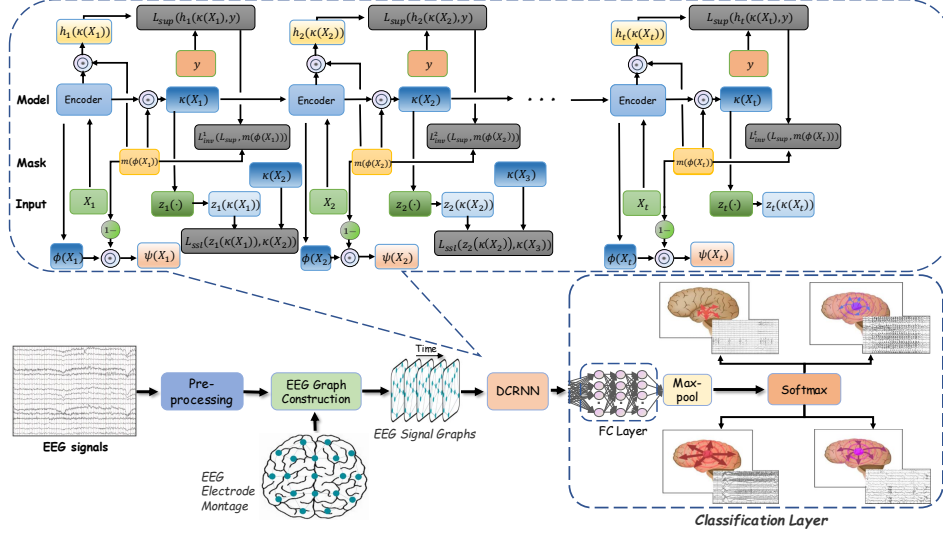


Figure 1: Overview of the proposed spatiotemporal invariant learning method.

84 An alternative approach, based on a directed correlation graph, particularly focuses on the dynamic  
 85 connectivity between different nodes. To evaluate the connectivity, only the weights that fall within  
 86 the most  $k$  similar neighbors (including self-edges) are preserved to ensure the sparsity of the graph.  
 87 The weight can be formulated as follows:

$$W_{ij} = \begin{cases} \text{Corr}(\mathbf{X}_{:,i,:}, \mathbf{X}_{:,j,:}) & \text{if } v_j \in \mathcal{C}_k(v_i) \\ 0 & \text{otherwise,} \end{cases}$$

88 where  $X_{:,i,:}$  and  $X_{:,j,:}$  denotes the preprocessed signals in  $v_i$  and  $v_j$ ,  $\text{Corr}(\cdot, \cdot)$  represents the pearson  
 89 correlation coefficient, and  $\mathcal{C}_k(v_i)$  referring to the most  $k$  similar neighbors of  $v_i$ .

90 **Diffusion Convolutional Recurrent Neural Network.** Previous works utilize the diffusion convolu-  
 91 tional recurrent neural network (DCRNN) to effectively capture the temporal and spatial dependencies  
 92 in EEG signals. To capture the temporal dependencies in EEG data, modified gated recurrent units  
 93 (GRUs) [5] are employed. For spatial dependency, diffusion convolution provides significant insights  
 94 into the influence exerted by each node on all others, and the extent of this kind of influence can be  
 95 quantified by applying a bidirectional random walk on the directed graph and calculating a  $K$ -step  
 96 diffusion convolution. The diffusion convolution is defined by:

$$X_{:,m*} \mathcal{G} f_{\theta} = \sum_{k=0}^{K-1} (\theta_{k,1} (D_O^{-1} W)^k + \theta_{k,2} (D_I^{-1} W^{\top})^k) X_{:,m}, \quad m \in \{1, \dots, M\},$$

97 where  $X$  is the preprocessed segment with  $N$  nodes and  $M$  features at timestamps  $t \in \{1, \dots, T\}$ ,  
 98  $\theta \in \mathbb{R}^{K \times 2}$  are the parameters of the filter, and  $D_O$  and  $D_I$  are the out-degree and in-degree diagonal  
 99 matrices of the graph. The transition matrices for the diffusion processes are  $D_O^{-1} W$  and  $D_I^{-1} W^{\top}$ .

100 For undirected graphs, the process converts to ChebNet spectral graph convolution [7], where  $X_{:,m}$   
 101 is filtered using Chebyshev polynomial bases. The spectral graph convolution can be expressed as

$$X_{:,m*} \mathcal{G} f_{\theta} = \Phi \left( \sum_{k=0}^{K-1} \theta_k \Lambda^k \right) \Phi^{\top} X_{:,m} = \sum_{k=0}^{K-1} \theta_k \mathbf{L}^k X_{:,m} = \sum_{k=0}^{K-1} \tilde{\theta}_k T_k(\tilde{L}) X_{:,m}, \quad m \in \{1, \dots, M\},$$

102 where  $\Phi$  and  $\Lambda$  are the eigenvector and eigenvalue matrix of the graph Laplacian  $\mathbf{L}$ .  $T_k(\tilde{L})$  is the  
 103  $k$ -th Chebyshev polynomial of the scaled Laplacian  $\tilde{L}$ , allowing for efficient computation without  
 104 explicit eigenvalue decomposition.

### 105 3 Methodology

106 In a cross-patient scenario, we propose the spatiotemporal invariant risk minimization (ST-IRM) loss,  
 107 making the prediction model achieves both (a) accurately predicting patient’s seizure type in each  
 108 patient group; (b) The variation of prediction between the different groups is small. Specifically,  
 109 for a timestamp  $t$ , we derive an invariant mask function  $m(\cdot)$  to separate the representations of the  
 110 raw EEG feature into two orthogonal components. We denote the representation of the raw EEG  
 111 feature as  $\phi(X_{:,:,t})$ . For simplification of notations, we use  $X_t$  instead of  $X_{:,:,t}$ . The representation  
 112 in the present paper is obtained by DRCNN. Through the invariant mask function  $m(\cdot)$ ,  $\phi(X_t)$  is  
 113 decomposed into an invariant representation  $\kappa(X_t) = m(\phi(X_t))$ , and the variant representation  
 114  $\psi(X_t) = (1 - m(\phi(X_t))) \odot \phi(X_t)$ , where  $m(X_t) \in [0, 1]^{N \times M}$ .

115 In time-series data, especially EEG data, there should be some correlation of the previous representa-  
 116 tions  $X_{t-1}$  with the current feature  $X_t$  [35]. Unlike the previous SSL approach that aims to learn  
 117 a model  $z_t(\cdot)$  to ensure  $z_{t-1}(X_{t-1}) \approx X_t$ , we claim that preserve the relation between the variant  
 118 parts,  $\psi(X_{t-1})$  and  $\psi(X_t)$  may not be helpful due to the spurious correlation. We expect only a good  
 119 prediction performance between the invariant representations. The proposed SSL loss is as below:

$$\mathcal{L}_{ssl} = \frac{1}{|nT|} \sum_{i=1}^n \sum_{t=1}^T \mathcal{L}(z_{t-1}(m(\phi(X_{t-1}^i))), m(\phi(X_t^i))),$$

120 where  $\mathcal{L}(\cdot, \cdot)$  is the loss function such as mean-square-error loss and  $X_t^i \in \mathbb{R}^{N \times M}$  is the preprocessed  
 121 signal for sample  $i$  at timestamp  $t$ . In addition, we want the information preserved by the mask  
 122 function can not only predict the next invariant representation but also can predict the final seizure  
 123 type, thus we use the following loss to provide the supervised signal for training the mask function:

$$\mathcal{L}_{sup} = \frac{1}{|n|} \sum_{i=1}^n \mathcal{L}(h_T(m(\phi(X_T^i))), y_i),$$

124 where  $h_T(\cdot)$  is the classification model and  $y_i$  is the ground truth label.

125 In addition, an ideal mask function  $m(\cdot)$  should be able to capture the invariant representation from  
 126 the raw EEG data. To address this, environments are created using the K-means clustering method to  
 127 separate the samples into groups, ensuring that samples within a group share similar characteristics,  
 128 while those in different groups exhibit distinct features. Thus, a classifier that performs consistently  
 129 across these environments would truly learn the invariant components and suffer the least from  
 130 spurious correlations. Assuming there is a total of  $G$  groups/environments, and the group indicator of  
 131 each sample is denoted by  $g_i$ . The supervised loss at timestamp  $t$  for the group  $g$  is given by

$$\mathcal{L}_{sup}^{g,t} = \frac{1}{\#\{i : g_i = g\}} \sum_{\{i:g_i=g\}} \mathcal{L}(h_t(m(\phi(X_t^i))), y_i),$$

132 where  $\#$  denotes the cardinal number of the set. It represents the supervised loss within the  $g$ -th group.  
 133 Combining the group-based supervised loss, the overall invariant risk loss at timestamp  $t$  is composed  
 134 of two major terms:

$$\mathcal{L}_{inv}^t = \mathbb{E}_{g \in \mathcal{G}} \mathcal{L}_{sup}^{g,t} + \lambda \|\text{Var}_{g \in \mathcal{G}} (\nabla_{\Theta^m} \mathcal{L}_{sup}^{g,t} \odot m(\phi(X_t)))\|^2,$$

135 where  $\Theta^m$  is the parameter of the mask function, and  $\lambda$  is the hyper parameter for tuning. The  
 136 previous term can be naively computed by  $\frac{1}{n} \sum_{g \in \mathcal{G}} \mathcal{L}_{sup}^{g,t}$ , suggesting the overall supervised loss at  
 137 timestamp  $t$ ; while the second term penalizes the classifier to perform consistently across groups. The  
 138 variance depicts the variation across the environments: the lower the variance is, the more consistent  
 139 performance the classifier obtains, thus, the better invariant presentation the classifier has learned  
 140 with. In the second term, we multiply the gradient with the mask function for scaling. For further  
 141 incorporating the spatiotemporal information, because the more information being observed, the  
 142 more accurate classification should be, we propose the weight decay loss below:

$$\mathcal{L}_{inv} = \sum_{t=1}^T w^{T-t} \mathcal{L}_{inv}^t,$$

143 where  $w \in (0, 1)$  is the weight decay rate, which is a hyper-parameter for tuning. The final proposed  
 144 ST-IRM loss is:

$$\mathcal{L}_{ST-IRM} = \mathcal{L}_{ssl} + \alpha \mathcal{L}_{sup} + \beta \mathcal{L}_{inv},$$

145 where  $\alpha$  and  $\beta$  are the hyper-parameters. An overview of the proposed method is given in Figure 1.

Table 1: Performance comparison of different methods under 12-second and 60-second scenario.

Method	12-s			60-s		
	F1	Recall	Precision	F1	Recall	Precision
CNN-LSTM	0.596 ± 0.035	0.654 ± 0.030	0.647 ± 0.036	0.623 ± 0.028	0.661 ± 0.030	0.647 ± 0.036
LSTM	0.690 ± 0.043	0.724 ± 0.033	0.725 ± 0.041	0.692 ± 0.011	0.718 ± 0.007	0.717 ± 0.017
Dense-CNN	0.657 ± 0.069	0.690 ± 0.053	0.694 ± 0.049	0.653 ± 0.085	0.704 ± 0.057	0.659 ± 0.118
MSTGCN	0.670 ± 0.031	0.719 ± 0.023	0.734 ± 0.029	0.647 ± 0.046	0.696 ± 0.027	0.694 ± 0.030
NeuroGNN	0.647 ± 0.040	0.710 ± 0.024	0.744 ± 0.030	0.698 ± 0.044	0.733 ± 0.042	0.714 ± 0.056
Corr-DCRNN	0.729 ± 0.058	0.756 ± 0.041	0.752 ± 0.047	0.672 ± 0.038	0.712 ± 0.021	0.705 ± 0.029
Dist-DCRNN	0.713 ± 0.044	0.735 ± 0.043	0.734 ± 0.045	0.695 ± 0.028	0.735 ± 0.013	0.738 ± 0.021
PANN-DCRNN	0.728 ± 0.052	0.753 ± 0.042	0.755 ± 0.041	0.684 ± 0.023	0.717 ± 0.016	0.720 ± 0.024
ST-InvDCRNN(ours)	<b>0.748 ± 0.038</b>	<b>0.772 ± 0.028</b>	<b>0.764 ± 0.043</b>	<b>0.713 ± 0.043</b>	<b>0.741 ± 0.024</b>	<b>0.742 ± 0.037</b>

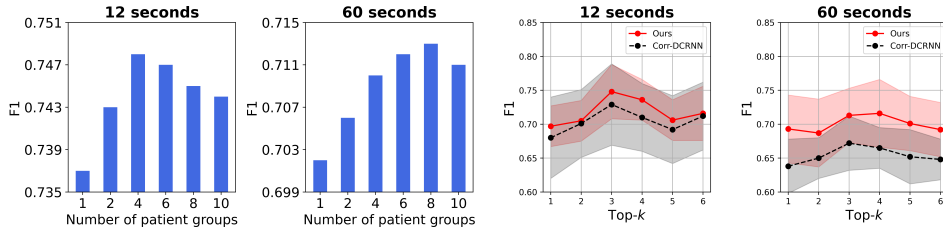


Figure 2: F1 under different numbers of patient groups (the two subfigures on the left) and different values of hyper-parameter top- $k$  to control the graph sparsity (the two subfigures on the right).

## 146 4 Experiments

### 147 4.1 Experimental Settings

148 **Datasets.** Following previous works [20, 31, 36], we utilized the Temple University Hospital EEG  
 149 Seizure Corpus (TUSZ) dataset, which is the largest public dataset for our experiments. Specifically,  
 150 we use the version v1.5.2 of the TUSZ dataset. The TUSZ dataset contains 5,612 EEG signals, and  
 151 3,050 annotated seizure events from over 300 patients, covering eight seizure types. The EEG signal  
 152 was recorded using 19 electrodes from the standard 10-20 system [14].

153 **Data preprocessing and Experiment Details.** Following the preprocessing approach of Tang et al.  
 154 [35], we transform the raw EEG signals into the frequency domain, as seizures are associated with  
 155 brain electrical activity in specific frequency bands [37]. Following prior methodologies [2, 3],  
 156 EEG recordings were resampled to 200Hz and segmented into non-overlapped 60-second windows  
 157 (clips), and only clips that contain a single type of seizure are considered. If a seizure event ends and  
 158 another begins within the same clip, it is truncated and zero-padded to preserve a 60-second duration.  
 159 Each 60-second clip is then segmented into 1-second intervals. The Fast Fourier Transform (FFT)  
 160 algorithm is applied to each segment to obtain the logarithmic amplitudes of non-negative frequency  
 161 components, as is outlined in Tang et al. [35]. Consequently, each 60-second clip is transformed  
 162 into a sequence of 60 log-amplitude vectors. In addition, following recent studies on seizure type  
 163 classification [2, 3, 35], we use weighted F1-score as the main evaluation metric with precision  
 164 and recall as well to measure the classification performance. See Appendix B for more experiment  
 165 protocols and details.

166 **Baselines.** We compare our proposed method with CNN-based method: **DenseCNN** [29], RNN-  
 167 based method: **LSTM** [13], and hybrid approach that combine CNN and RNN: **CNN-LSTM** [2].  
 168 We also compared our method with GNN-based methods: **MSTGCN** [16], **Dist-DCRNN** [35],  
 169 **Corr-DCRNN** [35], **NeuroGNN** [11], and **PANN** [44].

### 170 4.2 Performance Analysis

171 Table 1 shows the performance of our method compared with various baseline methods, evaluating  
 172 with three metrics, i.e., weighted F1, Recall, and Precision scores. First, DCRNN-based models  
 173 achieve competitive performance among all baselines. Second, our method significantly outperforms  
 174 the baselines under both scenarios with 12-second and 60-second clip windows. Note that we adopt  
 175 DCRNN as a backbone in the experiment, which is shown in Figure 1, and the superior against  
 176 DCRNN-based methods demonstrates the effectiveness of our invariant learning framework.

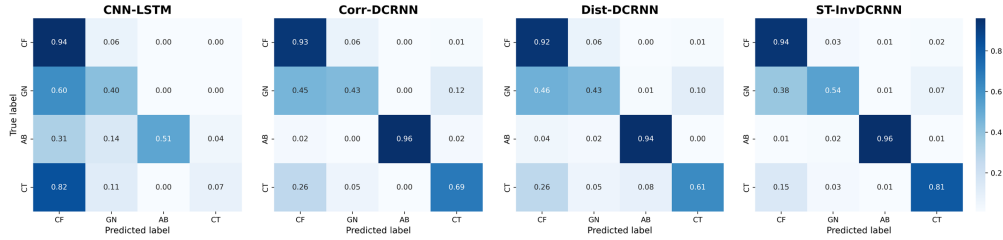


Figure 3: Confusion matrices for four classes of seizures.

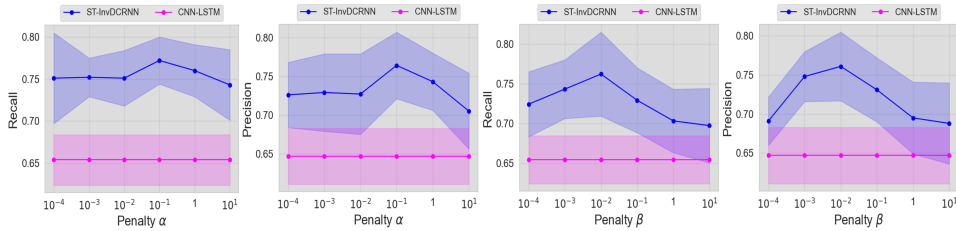


Figure 4: 12-second Performance under different penalty weights.

### 177 4.3 In-Depth Analysis

178 To comprehensively evaluate the proposed invariant learning method, we conduct four in-depth  
 179 analyses on the number of patient groups, the value of hyper-parameter top- $k$ , the classification  
 180 confusion matrix, and the 12-second performance under different penalty weights respectively. Note  
 181 that the patients are clustered into groups according to their EEG recordings, and the two subfigures  
 182 on the left of Figure 2 show that different numbers of groups result in varying performance. In  
 183 the scenario of 12-second clip windows, the best choice for group number is 4, while in the 60-  
 184 second case, the best value is 8. Our method outperforms Corr-DCRNN with top- $k$  ranging from  
 185 1 to 6, and the highest F1 is achieved when top- $k$  is around 3 for both scenarios. In addition, we  
 186 provide the results of the recall metric and the confusion metrics in Appendix B. Figure 3 shows the  
 187 confusion matrices for four seizure classification models, highlighting the superior performance of  
 188 our method. The ST-InvDCRNN reduces misclassifications and confusion between seizure types,  
 189 notably 0.81 for the CT class and 0.54 for GN seizures, outperforming baseline models. Figure 4  
 190 compares ST-InvDCRNN and CNN-LSTM performance across different penalty parameters ( $\alpha$  and  
 191  $\beta$ ) for recall and precision. ST-InvDCRNN consistently outperforms CNN-LSTM, especially at  
 192 intermediate penalty values. For Penalty  $\alpha$ , ST-InvDCRNN peaks at  $\alpha = 10^{-1}$ , achieving 0.772  
 193 recall score and 0.764 precision score, while CNN-LSTM shows lower scores. Similarly, for Penalty  
 194  $\beta$ , ST-InvDCRNN reaches its best performance at  $\beta = 10^{-1}$ , with 0.762 recall score and 0.761  
 195 precision score. Overall, ST-InvDCRNN delivers better classification results.

## 196 5 Conclusion

197 Epilepsy remains a significant global health challenge, with traditional EEG-based diagnostic methods  
 198 posing limitations due to their reliance on clinician review. With the recent advancement of deep  
 199 learning, techniques such as CNNs, RNNs, and GNNs are proposed to automatically classify the  
 200 seizure type. However, existing methods often lack cross-patient robustness and guarantee, which  
 201 is very common in practice. In addition, most of the methods addressing the cross-patient problem  
 202 ignore the spatiotemporal information. To bridge this gap, we propose a spatiotemporal invariant  
 203 risk minimization approach that addresses these challenges by adopting self-supervised learning and  
 204 capturing time-varying invariant features. Experimental results on the largest public dataset verify  
 205 the effectiveness of our approach, demonstrating its potential to advance epilepsy diagnosis in the  
 206 cross-patient scenario. One of the possible limitations is to investigate a more efficient way to learn  
 207 the model parameters and reduce the complexity while maintaining the classification performance.

208 **References**

- 209 [1] Arshia Afzal, Grigorios Chrysos, Volkan Cevher, and Mahsa Shoaran. Rest: Efficient and accelerated eeg  
210 seizure analysis through residual state updates. *arXiv preprint arXiv:2406.16906*, 2024.
- 211 [2] David Ahmedt-Aristizabal, Tharindu Fernando, Simon Denman, Lars Petersson, Matthew J Aburn, and  
212 Clinton Fookes. Neural memory networks for seizure type classification. In *2020 42nd Annual International  
213 Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 569–575. IEEE, 2020.
- 214 [3] Umar Asif, Subhrajit Roy, Jianbin Tang, and Stefan Harrer. SeizureNet: Multi-spectral deep feature  
215 learning for seizure type classification. In *Machine Learning in Clinical Neuroimaging and Radiogenomics  
216 in Neuro-oncology: Third International Workshop, MLCN 2020, and Second International Workshop,  
217 RNO-AI 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4–8, 2020, Proceedings 3*,  
218 pages 77–87. Springer, 2020.
- 219 [4] Ettore Beghi. The epidemiology of epilepsy. *Neuroepidemiology*, 54(2):185–191, 2020.
- 220 [5] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of  
221 neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop  
222 on Syntax, Semantics and Structure in Statistical Translation*, page 103. Association for Computational  
223 Linguistics, 2014.
- 224 [6] Alexander Craik, Yongtian He, and Jose L Contreras-Vidal. Deep learning for electroencephalogram (eeg)  
225 classification tasks: a review. *Journal of neural engineering*, 16(3):031001, 2019.
- 226 [7] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs  
227 with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- 228 [8] Jessica Falco-Walter. Epilepsy—definition, classification, pathophysiology, and epidemiology. In *Seminars  
229 in neurology*, volume 40, pages 617–623. Thieme Medical Publishers, Inc., 2020.
- 230 [9] Robert S Fisher, J Helen Cross, Jacqueline A French, Norimichi Higurashi, Edouard Hirsch, Floor E Jansen,  
231 Lieven Lagae, Solomon L Moshé, Jukka Peltola, Eliane Roulet Perez, et al. Operational classification  
232 of seizure types by the international league against epilepsy: Position paper of the ilae commission for  
233 classification and terminology. *Epilepsia*, 58(4):522–530, 2017.
- 234 [10] Arash Hajisafi, Haowen Lin, Sina Shaham, Haoji Hu, Maria Despoina Siampou, Yao-Yi Chiang, and  
235 Cyrus Shahabi. Learning dynamic graphs from all contextual information for accurate point-of-interest  
236 visit forecasting. In *Proceedings of the 31st ACM International Conference on Advances in Geographic  
237 Information Systems*, pages 1–12, 2023.
- 238 [11] Arash Hajisafi, Haowen Lin, Yao-Yi Chiang, and Cyrus Shahabi. Dynamic gnns for precise seizure  
239 detection and classification from eeg data. In *Pacific-Asia Conference on Knowledge Discovery and Data  
240 Mining*, pages 207–220. Springer, 2024.
- 241 [12] Jiatong He, Jia Cui, Gaobo Zhang, Mingrui Xue, Dengyu Chu, and Yanna Zhao. Spatial–temporal seizure  
242 detection with graph attention network and bi-directional lstm architecture. *Biomedical Signal Processing  
243 and Control*, 78:103908, 2022.
- 244 [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,  
245 1997.
- 246 [14] Richard W Homan, John Herman, and Phillip Purdy. Cerebral location of international 10–20 system  
247 electrode placement. *Electroencephalography and Clinical Neurophysiology*, 66(4):376–382, 1987. ISSN  
248 0013-4694.
- 249 [15] Chengbin Huang, Weiting Chen, and Guitao Cao. Automatic epileptic seizure detection via attention-based  
250 cnn-birnn. In *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages  
251 660–663, 2019. doi: 10.1109/BIBM47256.2019.8983420.
- 252 [16] Ziyu Jia, Youfang Lin, Jing Wang, Xiaojun Ning, Yuanlai He, Ronghao Zhou, Yuhan Zhou, and H Lehman  
253 Li-wei. Multi-view spatial-temporal graph convolutional networks with domain generalization for sleep  
254 stage classification. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 29:1977–1986,  
255 2021.
- 256 [17] Yizhang Jiang, Dongrui Wu, Zhaohong Deng, Pengjiang Qian, Jun Wang, Guanjin Wang, Fu-Lai Chung,  
257 Kup-Sze Choi, and Shitong Wang. Seizure classification from eeg signals using transfer learning, semi-  
258 supervised learning and tsf fuzzy system. *IEEE Transactions on Neural Systems and Rehabilitation  
259 Engineering*, 25(12):2270–2284, 2017.

- 260 [18] Dominik Klepl, Min Wu, and Fei He. Graph neural network-based eeg classification: A survey. *IEEE*  
261 *Transactions on Neural Systems and Rehabilitation Engineering*, 2024.
- 262 [19] Alicia Guadalupe Lazcano-Herrera, Rita Q Fuentes-Aguilar, and Mariel Alfaro-Ponce. Eeg motor/imagery  
263 signal classification comparative using machine learning algorithms. In *2021 18th International Conference*  
264 *on Electrical Engineering, Computing Science and Automatic Control (CCE)*, pages 1–6. IEEE, 2021.
- 265 [20] Yang Li, Yu Liu, Wei-Gang Cui, Yu-Zhu Guo, Hui Huang, and Zhong-Yi Hu. Epileptic seizure detection in  
266 eeg signals using a unified temporal-spectral squeeze-and-excitation network. *IEEE Transactions on Neural*  
267 *Systems and Rehabilitation Engineering*, 28(4):782–794, 2020. doi: 10.1109/TNSRE.2020.2973434.
- 268 [21] Jiashuo Liu, Zheyang Shen, Peng Cui, Linjun Zhou, Kun Kuang, Bo Li, and Yishi Lin. Stable adversarial  
269 learning under distributional shifts. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
270 volume 35, pages 8662–8670, 2021.
- 271 [22] Jiashuo Liu, Zheyang Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards  
272 out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.
- 273 [23] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint*  
274 *arXiv:1608.03983*, 2016.
- 275 [24] X Lou, X Li, H Meng, J Hu, M Xu, Y Zhao, J Yang, and Z Li. Eeg-dbnnet: A dual-branch network for  
276 temporal-spectral decoding in motor-imagery brain-computer interfaces. 2024.
- 277 [25] Yahong Ma, Zhenhao Huang, Jianyun Su, Hangyu Shi, Dong Wang, Shanshan Jia, and Weisu Li. A multi-  
278 channel feature fusion cnn-bi-lstm epilepsy eeg classification and prediction model based on attention  
279 mechanism. *IEEE Access*, 11:62855–62864, 2023.
- 280 [26] Bijan Mazaheri, Atalanti Mastakouri, Dominik Janzing, and Michaela Hardt. Causal information splitting:  
281 Engineering proxy features for robustness to distribution shifts. In *Uncertainty in Artificial Intelligence*,  
282 pages 1401–1411. PMLR, 2023.
- 283 [27] Michael Oberst, Nikolaj Thams, Jonas Peters, and David Sontag. Regularizing towards causal invariance:  
284 Linear models with proxies. In *International Conference on Machine Learning*, pages 8260–8270. PMLR,  
285 2021.
- 286 [28] Advait U Parulekar, Karthikeyan Shanmugam, and Sanjay Shakkottai. Pac generalization via invariant  
287 representations. In *International Conference on Machine Learning*, pages 27378–27400. PMLR, 2023.
- 288 [29] Khaled Saab, Jared Dunnmon, Christopher Ré, Daniel Rubin, and Christopher Lee-Messer. Weak super-  
289 vision as an efficient approach for automated seizure detection in electroencephalography. *NPJ digital*  
290 *medicine*, 3(1):59, 2020.
- 291 [30] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural  
292 networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint*  
293 *arXiv:1911.08731*, 2019.
- 294 [31] Rijad Sarić, Dejan Jokić, Nejra Beganović, Lejla Gurbeta Pokvić, and Almir Badnjević. Fpga-based  
295 real-time epileptic seizure classification using artificial neural network. *Biomedical Signal Processing and*  
296 *Control*, 62:102106, 2020. ISSN 1746-8094.
- 297 [32] MNAH Sha’ Abani, N Fuad, Norezmi Jamal, and MF Ismail. knn and svm classification for eeg: a review.  
298 In *InECCE2019: Proceedings of the 5th International Conference on Electrical, Control & Computer*  
299 *Engineering, Kuantan, Pahang, Malaysia, 29th July 2019*, pages 555–565. Springer, 2020.
- 300 [33] Afshin Shoeibi, Marjane Khodatars, Navid Ghassemi, Mahboobeh Jafari, Parisa Moridian, Roohallah  
301 Alizadehsani, Maryam Panahiazar, Fahime Khozeimeh, Assef Zare, Hossein Hosseini-Nejad, et al. Epilep-  
302 tic seizures detection using deep learning techniques: a review. *International journal of environmental*  
303 *research and public health*, 18(11):5780, 2021.
- 304 [34] Supriya Supriya, Siuly Siuly, Hua Wang, and Yanchun Zhang. Epilepsy detection from eeg using complex  
305 network techniques: A review. *IEEE Reviews in Biomedical Engineering*, 16:292–306, 2021.
- 306 [35] Siyi Tang, Jared Dunnmon, Khaled Kamal Saab, Xuan Zhang, Qianying Huang, Florian Dubost, Daniel  
307 Rubin, and Christopher Lee-Messer. Self-supervised graph neural networks for improved electroencephalo-  
308 graphic seizure analysis. In *International Conference on Learning Representations*, 2022.



- 309 [36] Punrawish Thuwajit, Phurin Rangpong, Phattarapong Sawangjai, Phairot Autthasan, Rattanaphon  
310 Chaisaen, Nannapas Banluesombatkul, Puttaranun Boonchit, Nattasate Tatsaringkansakul, Thapanun  
311 Sudhawiyangkul, and Theerawit Wilaiprasitporn. Eegwavenet: Multiscale cnn-based spatiotemporal  
312 feature extraction for eeg seizure detection. *IEEE Transactions on Industrial Informatics*, 18(8):5547–5557,  
313 2022. doi: 10.1109/TII.2021.3133307.
- 314 [37] Alexandros T Tzallas, Markos G Tsipouras, and Dimitrios I Fotiadis. Epileptic seizure detection in  
315 eegs using time–frequency analysis. *IEEE transactions on information technology in biomedicine*, 13(5):  
316 703–710, 2009.
- 317 [38] Haiyang Yang, Shixiang Tang, Meilin Chen, Yizhou Wang, Feng Zhu, Lei Bai, Rui Zhao, and Wanli  
318 Ouyang. Domain invariant masked autoencoders for self-supervised learning from multi-domains. In  
319 *European Conference on Computer Vision*, pages 151–168. Springer, 2022.
- 320 [39] Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. Causalvae: Dis-  
321 entangled representation learning via neural structural causal models. In *Proceedings of the IEEE/CVF*  
322 *conference on computer vision and pattern recognition*, pages 9593–9602, 2021.
- 323 [40] Zhizhang Yuan, Daoze Zhang, YANG YANG, Junru Chen, and Yafeng Li. Ppi: Pretraining brain signal  
324 model for patient-independent seizure detection. In A. Oh, T. Naumann, A. Globerson, K. Saenko,  
325 M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages  
326 69586–69604. Curran Associates, Inc., 2023.
- 327 [41] Zhang Yutian, Huang Shan, Zhang Jianing, and Fan Ci’en. Design and implementation of an emotion  
328 analysis system based on eeg signals. *arXiv preprint arXiv:2405.16121*, 2024.
- 329 [42] Xiang Zhang, Lina Yao, Manqing Dong, Zhe Liu, Yu Zhang, and Yong Li. Adversarial representation  
330 learning for robust patient-independent epileptic seizure detection. *IEEE journal of biomedical and health*  
331 *informatics*, 24(10):2852–2859, 2020.
- 332 [43] Zeyang Zhang, Xin Wang, Ziwei Zhang, Haoyang Li, Zhou Qin, and Wenwu Zhu. Dynamic graph neural  
333 networks under spatio-temporal distribution shift. *Advances in neural information processing systems*, 35:  
334 6074–6089, 2022.
- 335 [44] Zongpeng Zhang, Taoyun Ji, Mingqing Xiao, Wen Wang, Guojing Yu, Tong Lin, Yuwu Jiang, Xiaohua  
336 Zhou, and Zhouchen Lin. Cross-patient automatic epileptic seizure detection using patient-adversarial  
337 neural networks with spatio-temporal eeg augmentation. *Biomedical Signal Processing and Control*, 89:  
338 105664, 2024.

## 339 Appendix

### 340 A Experimental Details

341 Following previous works, we divide the clips and patients of the TUSZ dataset into training,  
342 validation, and test sets. The number of EEG clips is 1,925, 450, and 521 for the three sets respectively,  
343 while the number of patients is 179, 22, and 34. Note that the patient sets are disjoint for training,  
344 validation, and test sets to study the cross-patient seizure classification robustness.

345 We tune the following hyper-parameters on the validation set.

- 346 •  $lr_{init} \in [1e - 5, 5e - 3]$ , the initial learning rate;
- 347 •  $top-k \in \{1, 2, 3, 4, 5, 6\}$ , the number of neighbors included in the correlation graphs for  
348 each node;
- 349 •  $K \in \{2, 3, 4\}$ , the maximum diffusion step;
- 350 •  $d \in [0, 0.7]$ , the dropout probability in the prediction networks.
- 351 •  $e \in [20, 40, 60, 80, 100]$ , the number of training epochs.

352 During the training, each batch has 40 EEG clips and the cosine annealing learning rate scheduler [23]  
353 is adopted. Our experiments are conducted on a computing platform of NVIDIA GeForce RTX 3090  
354 and Intel(R) Xeon(R) Gold 6248R CPU @ 3.00GHz.

355 **B Additional Evaluation Results**

356 Figure 5 shows the weighted F1 and the Recall scores to evaluate the performance of our method  
 357 under different number of patient groups, for both 12-second and 60-second clip windows. We can  
 358 observe that as the number of patient groups increases, the Recall-score has a similar pattern as the  
 359 weighted F1-score, achieving the highest value at 4 for the 12-second case and 8 for the 60-second  
 case.

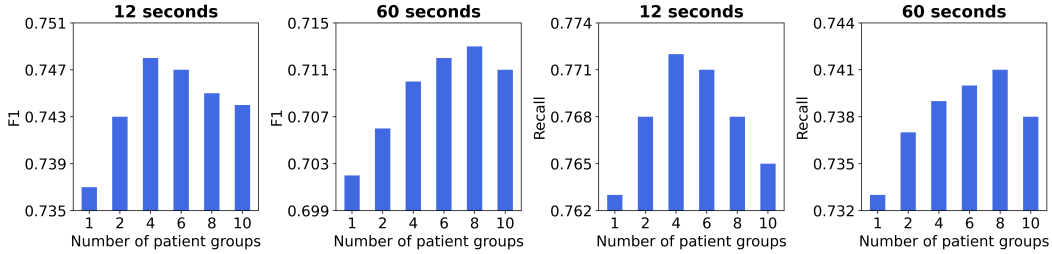


Figure 5: Performance under different numbers of patient groups.

360

361 Figure 6 shows the weighted F1 and the Recall scores to compare the performance of our method  
 362 with Corr-DCRNN under different top- $k$  values, for both 12-second and 60-second clip windows. As  
 363 the value of top- $k$  ranges from 1 to 6, the trend for both weighted F1 and Recall scores is increasing  
 364 until a peak at around 3, followed by a slight decrease.

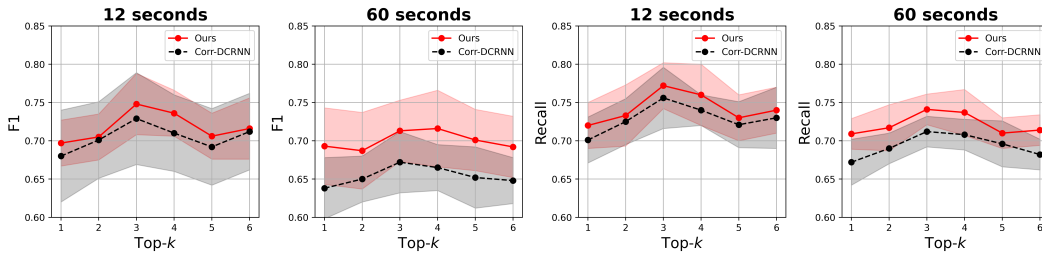


Figure 6: Performance under different values of top- $k$ .