

Con4m: UNLEASHING THE POWER OF CONSISTENCY AND CONTEXT IN CLASSIFICATION FOR BLURRED-SEGMENTED TIME SERIES

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Blurred-Segmented Time Series (BST) has emerged as a prevalent form of time
 2 series data in various practical applications, presenting unique challenges for the
 3 Time Series Classification (TSC) task. The BST data is segmented into continuous
 4 states with inherently blurred transitions. These transitions lead to inconsistency
 5 in annotations among different individuals due to experiential differences, thereby
 6 hampering model training and validation. However, existing TSC methods often
 7 fail to recognize label inconsistency and contextual dependencies between consec-
 8 utive classified samples. In this work, we first theoretically clarify the connotation
 9 of valuable contextual information. Based on these insights, we incorporate prior
 10 knowledge of BST data at both the data and class levels into our model design to
 11 capture effective contextual information. Furthermore, we propose a label consistency
 12 training framework to harmonize inconsistent labels. Extensive experiments
 13 on two public and one private BST data fully validate the effectiveness of our pro-
 14 posed approach, *Con4m*, in handling the TSC task on BST data.

15 1 INTRODUCTION

16 Time series classification (TSC) has been widely studied in the field of machine learning for many
 17 years (Middlehurst et al., 2023). With the rapid development of measurement technology re-
 18 cently, TSC has been extended to various applications in diverse practical domains, such as health-
 19 care (Rafiei et al., 2022; Chen et al., 2022), finance (Dezhkam et al., 2022; Liu & Cheng, 2023), and
 20 environmental monitoring (Yuan et al., 2022; Tian et al., 2023). TSC often involves in classifying
 21 time series samples into predefined categories with labels and is usually based on the assumption of
 22 independence and identical distribution (*i.i.d.*) (Dempster et al., 2021; Zhao et al., 2023).

23 In practical applications, however, a large number of **Blurred-Segmented Time Series (BST)** data
 24 have emerged, which differ in fundamental ways from traditional TSC data: **(1) BST intrinsically
 25 records blurred transitions on the boundaries between different states.** For example, in terms
 26 of a person’s emotional state, the transition from sadness to happiness is ambiguous, with no clear
 27 boundaries. **(2) States last for a long duration, segmenting BST.** Take sleep data covering physi-
 28 ological signals of subjects overnight as an example, it shows alternations of different sleep stages,
 29 each of which stably lasts for a prolonged period.

30 The characteristics of BST pose new challenges for mainstream TSC models. **Firstly, the presence
 31 of blurred boundaries leads to inconsistent annotations.** In the case of raw BST data, manual an-
 32 notations usually determine the start and end points of a particular state. Especially in the healthcare
 33 domain, data is collected from different hospitals. Due to the lack of standardized quantification cri-
 34 teria, annotations from different doctors vary for their individual experiences. **In the TSC task, each
 35 type of states is assigned a unique label.** Therefore, the inconsistency in labeling across different
 36 data sources hampers model training. However, most existing TSC works model time series data by
 37 assuming noise-free labels, which significantly limits their performance on BST data.

38 **The continuous states and gradual transitions call for more coherent contextual prediction.** In
 39 the TSC task, BST data is divided into time segments corresponding to different states (labels) to
 40 be classified. There are natural temporal dependencies between consecutive segments, which not
 41 only exists at the data level but also manifests in the changes of labels. However, mainstream TSC

models (Middlehurst et al., 2023; Foumani et al., 2023) are often designed for publicly available datasets (Bagnall et al., 2018; Dau et al., 2019) based on *i.i.d.* samples, disregarding the inherent contextual dependencies between the samples in time series data. Although some time series models (Shao et al., 2022; Nie et al., 2023) take contextual information of the input data into consideration for predictions with patch-by-patch modeling, they fail to incorporate the class information of consecutive classified time segments so as to achieve coherent predictions for BST data.

To better model BST data, we first analyze how to enhance the relevance between input data and labels in classification tasks by introducing effective contextual information from an information-theoretic perspective. Subsequently, based on the theoretic insights, we incorporate contextual prior knowledge of BST data from both the data and label perspectives to improve the prediction ability of the model. Lastly, drawing inspiration from noisy label learning, we enable the model to progressively harmonize inconsistent labels during the learning process of classification. Consequently, we propose *Con4m* (pronounced **Conform**) - a label **C**onsistency training framework with effective **C**ontextual information, achieving **C**oherent predictions and **C**ontinuous representations for time series classification on BST data. Extensive experiments on two public and one private BST data demonstrate the superior performance of *Con4m*. In addition, we verify the *Con4m*'s ability to harmonize inconsistent labels by the label substitution experiment. A case study is also shown to give further insight into how *Con4m* works well for BST data.

Our contributions are as follows. **(1)** We are the first to emphasize the importance of BST data and systematically analyze and model it, which is critical for various practical applications. **(2)** We theoretically elucidate the valuable contextual information for the input data in the classification task. Combined with the theoretical insights, we propose a novel framework *Con4m* that can be effectively applied to the TSC task with BST data. **(3)** Extensive experiments fully highlight the superiority of *Con4m* for modeling BST data, shedding light on the era of personalized services when applications like precision medicine, physiological status monitoring and others will prevail.

2 VALUABLE CONTEXTS ENHANCE PREDICTIVE ABILITY

Intuitively, it is widely believed that the performance of models on the classification task can be enhanced by incorporating contextual information. But why does this conclusion hold? What kind of contextual information should be introduced? In this section, we aim to analyze this phenomenon from an information-theoretic perspective at a macro level.

Assuming that the random variables of the classified samples and their corresponding labels are denoted as x_t and y_t , respectively. \mathbb{A}_t represents the contextual sample set introduced for x_t . $x_{\mathbb{A}_t}$ denotes the random variable for the contextual sample set.

Proposition 1. *Introducing contextual information does not compromise the performance of a model for the classification task.*

Proof.

$$\mathbb{I}(y_t; x_t, x_{\mathbb{A}_t}) = \mathbb{I}(y_t; x_{\mathbb{A}_t} | x_t) + \mathbb{I}(y_t; x_t) \geq \mathbb{I}(y_t; x_t). \quad (1)$$

The inequality holds due to the non-negativity of conditional mutual information. **Mutual information measures the correlation between two variables. In the classification task, a higher correlation between samples and labels indicates that the samples are more easily distinguishable by the labels. Based on the assumption that a model can perfectly capture these correlations, a higher mutual information implies a higher upper bound on the model's performance in classifying samples.** \square

According to (1), the increase in $I(y_t; x_{\mathbb{A}_t} | x_t)$ determines the extent to which the upper bound of the model's performance improves. Hence, we employ Theorem 1 to elucidate the specific contextual sample set that can maximize the information gain $I(y_t; x_{\mathbb{A}_t} | x_t)$.

Theorem 1. *Introducing a contextual sample set that maximizes the predictive ability of labels yields the maximum information gain.*

Proof. Expanding $\mathbb{I}(y_t; x_{\mathbb{A}_t} | x_t)$, we have:

$$\mathbb{I}(y_t; x_{\mathbb{A}_t} | x_t) = \sum_{x_t} p(x_t) \sum_{x_{\mathbb{A}_t}} \sum_{y_t} p(y_t, x_{\mathbb{A}_t} | x_t) \log \frac{p(y_t, x_{\mathbb{A}_t} | x_t)}{p(y_t | x_t) p(x_{\mathbb{A}_t} | x_t)} \quad (2)$$

$$= \sum_{\mathbf{x}_t} p(\mathbf{x}_t) \sum_{\mathbf{x}_{\mathbb{A}_t}} \sum_{y_t} p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t}) p(\mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) \log \frac{p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t})}{p(y_t|\mathbf{x}_t)} \quad (3)$$

$$= \sum_{\mathbf{x}_t} p(\mathbf{x}_t) \sum_{\mathbf{x}_{\mathbb{A}_t}} p(\mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) D_{\text{KL}}(p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t})||p(y_t|\mathbf{x}_t)). \quad (4)$$

88 Given a fixed classification sample \mathbf{x}_t and the inherent distribution $p(y_t|\mathbf{x}_t)$ of the data, the
 89 KL divergence is a convex function that attains its minimum at $p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t}) = p(y_t|\mathbf{x}_t)$. As
 90 $p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t})$ approaches the boundary of the probability space, indicating a stronger predictive
 91 ability for y_t , the value of KL divergence increases. Due to the convexity of KL divergence,
 92 there exists a contextual sample set in the data that maximizes $D_{\text{KL}}(p(y_t|\mathbf{x}_t, \mathbf{x}_{\mathbb{A}_t})||p(y_t|\mathbf{x}_t))$. We
 93 denote this sample set as \mathbb{A}_t^* and the maximum KL divergence value as D_t^* . Additionally, we
 94 note that $\sum_{\mathbf{x}_{\mathbb{A}_t}} p(\mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) = 1$. Hence, we can obtain the upper bound for the information gain
 95 $\mathbb{I}(y_t; \mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) \leq \sum_{\mathbf{x}_t} p(\mathbf{x}_t) \sum_{\mathbf{x}_{\mathbb{A}_t}} p(\mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) D_t^* \leq \sum_{\mathbf{x}_t} p(\mathbf{x}_t) D_t^*$.

96 To achieve this upper bound, the model needs to introduce a contextual sample set \mathbb{A}_t^* for each
 97 sample that maximally enhances its label’s predictive ability. Moreover, the model needs to reach
 98 an optimal selection strategy distribution $p(\mathbf{x}_{\mathbb{A}_t^*}|\mathbf{x}_t) = 1, p(\mathbf{x}_{\mathbb{A}_t}|\mathbf{x}_t) = 0$ (for $\mathbb{A}_t \neq \mathbb{A}_t^*$). \square

99 According to Theorem 1, the model needs to find the optimal contextual sample set that enhances its
 100 predictive ability. In this paper, we utilize learnable weights to allow the model to adaptively select
 101 potential contextual samples. Through explicit supervised learning, the model can directly enhance
 102 its predictive ability in an end-to-end manner. On the other hand, benefiting from an information-
 103 theoretic perspective, $\mathbf{x}_{\mathbb{A}_t}$ not only includes the raw data of contextual samples but also incorporates
 104 their label information, which can be represented as $y_{\mathbb{A}_t}$. Therefore, we can introduce contextual
 105 information at both the data and class levels to enhance the model’s predictive ability.

106 3 THE *Con4m* METHOD

107 In this section, we introduce the details of *Con4m*. Based on the insights of Theorem 1, we introduce
 108 effective contextual information at both the data (Sec. 3.1) and class (Sec. 3.2) levels to enhance the
 109 predictive ability of *Con4m*. In Sec. 3.3, inspired by the idea of noisy label learning, we propose
 110 a label harmonization framework to achieve label consistency. Before delving into the details of
 111 *Con4m*, we first provide the formal definition of the time series classification task in our work.

112 **Definition 1.** *Given a time interval comprising of T consecutive time points, denoted as $s =$
 113 $\{s_1, s_2, \dots, s_T\}$, a w -length sliding window with stride length r is employed for segmentation. s
 114 is partitioned into L time segments, represented as $x = \{x_i = \{s_{(i-1) \times r + 1}, \dots, s_{(i-1) \times r + w}\} | i =$
 115 $1, \dots, L\}$. The model is tasked with predicting labels for each time segment (sample) x_i .*

116 3.1 CONTINUOUS CONTEXTUAL REPRESENTATION ENCODER

117 BST data exhibits temporal persistence for each class. By paying closer attention to and aggregating
 118 neighboring segments, the model can acquire temporally smoother representations of time segments.
 119 Smoother representations lead to smoother predictive probabilities. This benefits not only the pre-
 120 diction of consecutive time segments belonging to the same class with the same label but also aligns
 121 with the gradual nature of class transitions. Therefore, we introduce the Gaussian prior to allow for
 122 a more targeted selection of the contextual sample set \mathbb{A}_t to enhance the model’s predictive ability.

123 Self-attention in BERT (Devlin et al., 2019) has the ability to globally model sequences. How-
 124 ever, point-wise attention computations often fail to obtain smooth representations after aggregation.
 125 Therefore, similar to the Gaussian filter technique, we use the Gaussian kernel $\Phi(x, y|\sigma)$ as prior
 126 weights to aggregate neighbors to obtain smoother representations. Since the neighbors of boundary
 127 segments may belong to different classes, we allow each segment to learn its own scale parameter
 128 σ . Formally, as Figure 1 shows, the two-branch **Con-Attention** in the l -th layer is:

$$Q, K, V_t, \sigma, V_s = c^{l-1} W_Q^l, c^{l-1} W_K^l, c^{l-1} W_{V_t}^l, c^{l-1} W_\sigma^l, c^{l-1} W_{V_s}^l, \quad (5)$$

$$T^l = \text{SoftMax} \left(\frac{QK^T}{\sqrt{d}} \right), \quad (6)$$

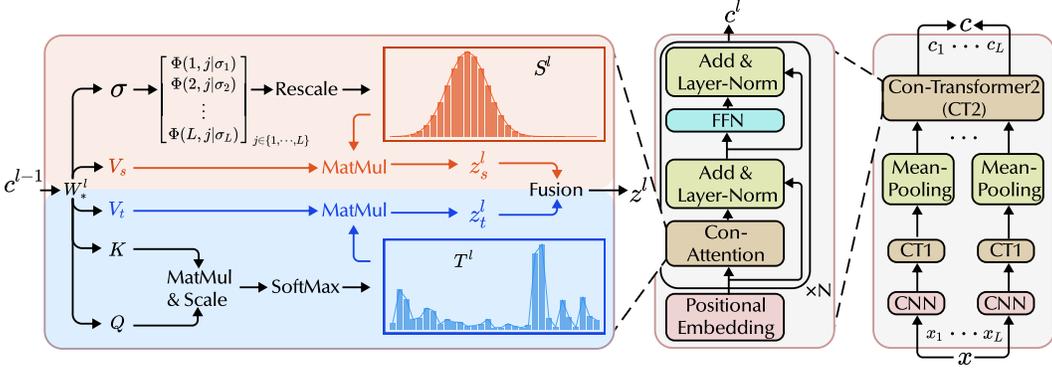


Figure 1: Overview of the encoder of *Con4m*. The leftmost part shows the details of Con-Attention. The right part of the figure shows the architecture of Con-Transformer and the whole encoder of *Con4m*.

$$S^l = \text{Rescale} \left(\left[\frac{1}{\sqrt{2\pi}\sigma_i} \exp \left(-\frac{|j-i|^2}{2\sigma_i^2} \right) \right]_{i,j \in \{1, \dots, L\}} \right), \quad (7)$$

$$z_t^l = T^l V_t, \quad z_s^l = S^l V_s, \quad (8)$$

129 where L is the number of the consecutive segments, d is the dimension of the hidden representations,
 130 $c^{l-1} \in \mathbb{R}^{L \times d}$ is the $l-1$ -th layer's hidden representations and $W_*^l \in \mathbb{R}^{d \times d}$ are all learnable matrices.
 131 $\text{Rescale}(\cdot)$ refers to row normalization by index i . Q , K and V vectors represent the query, key
 132 and value of the self-attention mechanism respectively. To distinguish between two computational
 133 branches, we use s/S to represent the branch based on Gaussian prior, and t/T to represent the
 134 branch based on vanilla self-attention. T^l and S^l are the aggregation weights of the two branches.

135 Then we use the conventional attention mechanism (Bahdanau et al., 2015) to adaptively fuse z_t^l and
 136 z_s^l . Finally, as illustrated in Figure 1, by stacking the multi-head version of Con-Attention layers,
 137 we construct Con-Transformer, which serves as the backbone of the continuous encoder of *Con4m*
 138 to obtain final representations c . During the practical implementation, we adopt the same approach
 139 proposed by Xu et al. (2022) for the computation of Gaussian kernel function.

140 3.2 CONTEXT-AWARE COHERENT CLASS PREDICTION

141 In the classification task of BST data, consecutive time segments not only provide context at the
 142 data level but also possess their own class information. For instance, in the case of human motion
 143 recognition, if an individual is walking at the beginning and end within a reasonable time range, it is
 144 highly likely that the intermediate states also corresponds to walking. Existing TSC models (Mid-
 145 dlehurst et al., 2023; Foumani et al., 2023) primarily focus on classifying independent segments,
 146 overlooking the temporal dependencies of the labels. **But our theoretic framework allows for the**
 147 **incorporation of contextual information at the class level into the model's design.**

148 **Neighbor Class Consistency Discrimination.** According to Theorem 1, we aim to identify a set
 149 of contextual samples that maximizes the model's predictive ability at the class level. Since directly
 150 optimizing the label aggregation is challenging, we adopt the approach of aggregating predictions
 151 of segments belonging to the same class. The idea is inspired by the observation that for graph
 152 neural networks based on the homophily assumption, aggregating neighbor information belonging
 153 to the same class can improve predictive performance (McPherson et al., 2001; Zhu et al., 2020).
 154 Therefore, we train a discriminator to determine whether two segments belong to the same class.
 155 The model then selects a contextual sample set based on the discriminator's predictions. As the left
 156 part of Figure 2 shows, we formalize this process as the following equations:

$$\hat{p} = \text{SoftMax}(\text{MLP}_1(c)), \quad Q, K, V = c, c, \hat{p}, \quad (9)$$

$$\hat{R} = \text{SoftMax} \left(\left[\text{MLP}_2(Q_i \| K_j) \right]_{i,j \in \{1, \dots, L\}} \right), \quad (10)$$

$$\tilde{p} = \hat{R}_{:,1} V, \quad (11)$$

157 where $\hat{R} \in \mathbb{R}^{L \times L \times 2}$ is the probability of whether two neighbor segments belong to the same class
 158 and $(\cdot \parallel \cdot)$ denotes tensor concatenation. We define the two losses for the model training as:

$$\ell_1 = \text{CrossEntropy}(\hat{p}, y), \quad \ell_2 = \text{CrossEntropy}(\hat{R}, \tilde{Y}), \quad (12)$$

159 where $\tilde{Y} = [\mathbf{1}_{y_i=y_j}]_{i,j \in \{1, \dots, L\}}$. Given that ℓ_1 and ℓ_2 are of the same magnitude, we equally sum
 160 them as the final loss.

161 **Prediction Behavior Constraint.** Although we incorporate the contextual class information, we
 162 still cannot guarantee the overall predictive behavior of consecutive segments. For the BST data,
 163 within a suitably chosen time interval, **the majority of** consecutive time segments span at most two
 164 classes. Therefore, the predictions in the intervals should exhibit a constrained monotonicity.

165 As shown in Figure 2, for each class in prediction results, there are only four prediction behaviors
 166 for consecutive time segments, namely *high confidence*, *low confidence*, *confidence decreasing*, and
 167 *confidence increasing*. To constrain the behavior, we use function fitting to integrate \tilde{p} . Considering
 168 the wide applicability, we opt for the hyperbolic tangent function (*i.e.*, Tanh) as our basis. Formally,
 169 we introduce four tunable parameters to exactly fit the monotonicity as:

$$\bar{p} = \text{Tanh}(x|a, k, b, h) = a \times \text{Tanh}(k \times (x + b)) + h, \quad (13)$$

170 where parameter a constrains the range of the function’s values, k controls the slope of the transition
 171 of the function, b and h adjust the symmetry center of the function, and x is the given free vector in
 172 the x-coordinate. We use the MSE loss to fit the contextual predictions \tilde{p} as follows:

$$\ell_3 = \|\text{Tanh}(x|a, k, b, h) - \tilde{p}\|^2. \quad (14)$$

173 It deserves to emphasize that \tilde{p} in this fit has no gradient and therefore does not affect the parameters
 174 of the encoder. Please see Appendix B for more fitting details.

175 After function fitting, we obtain independent predictions \hat{p} for each segment and constrained predictions
 176 \bar{p} that leverage the contextual class information. For the inference stage, we use the average
 177 of them as the final coherent predictions, *i.e.*, $\hat{y} = \arg \max (\hat{p} + \bar{p})/2$. Next, we demonstrate how
 178 these predictions are combined during the training phase to achieve harmonized labels.

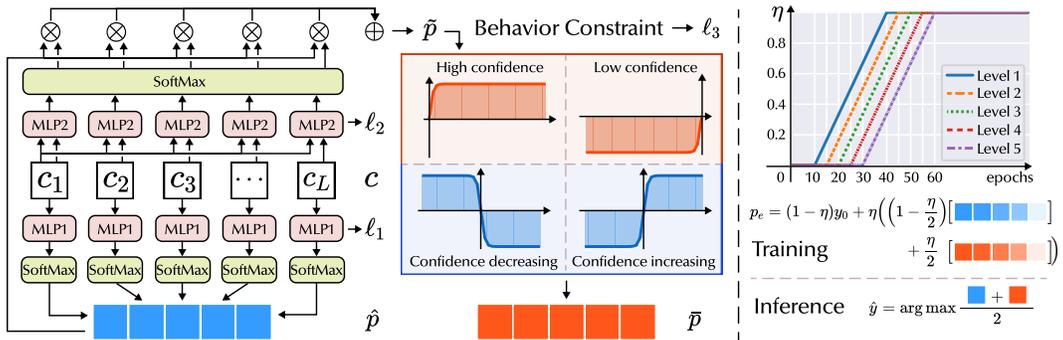


Figure 2: Overview of context-aware coherent class prediction and consistent label training framework in *Con4m*. The left part describes the neighbor class consistency discrimination task and the prediction behavior constraint. The rightmost part presents the training and inference details for label harmonization.

179 3.3 CONSISTENT LABEL TRAINING FRAMEWORK

180 Due to inherent blurred boundary, the annotation of BST data often lacks quantitative criteria, result-
 181 ing in experiential differences among individuals. Such discrepancies are detrimental to models and
 182 we propose a training framework to enable the model to adaptively harmonize inconsistent labels.

183 **Learning from easy to hard.** We are based on a fact that although people may have differences for
 184 the blurred transitions between states, they tend to reach an agreement on the most significant core
 185 part of the states. In other words, the empirical differences become more apparent when approaching
 186 the transitions. Therefore, we adopt curriculum learning techniques to help the model learn sam-
 187 ples from the easy (core) to the hard (transition) part. Formally (see the diagram in Figure 6(a) in

188 Appendix), for a continuous K -length state, we divide it into $N_l = 5$ equally sized levels as follows:

$$\left(\left[\lceil (N_l - 1) \frac{K}{2N_l} \rceil, \lfloor (N_l + 1) \frac{K}{2N_l} \rfloor \right]; \dots; \left[1, \lceil \frac{K}{2N_l} \rceil \right] \cup \left(\lfloor (2N_l - 1) \frac{K}{2N_l} \rfloor, K \right]. \quad (15)$$

189 Then we sample the same number of time intervals from each level. The higher the level, the more
190 apparent the inconsistency. Therefore, as Figure 2 shows, during the training stage, *Con4m* learns
191 the corresponding intervals in order from low to high levels, with a lag gap of $E_g = 5$ epochs.

192 **Harmonizing inconsistent labels.** Inspired by the idea of noisy label learning, we gradually
193 change the original labels to harmonize the inconsistency. The model preferentially changes the
194 labels of the core segments that are easier to reach a consensus, which can avoid overfitting of
195 uncertain labels. Moreover, the model will consider both the independent and contextual predictions
196 to robustly change inconsistent labels. Specifically, given the initial label y_0 , we update the labels
197 $y_e = \arg \max p_e$ for the e -th epoch, where p_e is obtained as follows:

$$\omega(e, 5) = \text{Rescale} \left(\left[\exp \left((e - m) / 2 \right) \right]_{m \in \{0, \dots, 4\}} \right), \quad (16)$$

$$\hat{p}_e^5 = \omega(e, 5) \cdot [\hat{p}_{e-m}]_{m \in \{0, \dots, 4\}}, \quad \bar{p}_e^5 = \omega(e, 5) \cdot [\bar{p}_{e-m}]_{m \in \{0, \dots, 4\}}, \quad (17)$$

$$p_e = (1 - \eta) y_0 + \eta \left(\left(1 - \frac{\eta}{2} \right) \hat{p}_e^5 + \frac{\eta}{2} \bar{p}_e^5 \right), \quad (18)$$

198 where \hat{p}_{e-m} and \bar{p}_{e-m} are the independent and contextual predictions in the $e - m$ -th epoch respec-
199 tively and \cdot denotes the dot product. $\omega(e, 5)$ is the exponentially averaged weight vector to aggregate
200 the predictions of the last 5 epochs to achieve more robust label update. The dynamic weighting fac-
201 tor, η , is used to adjust the degree of label update. As Figure 2 shows, η linearly increases from 0
202 to 1 with E_η epochs, gradually weakening the influence of the original labels. Besides, in the initial
203 training stage, the model tends to improve independent predictions. As the accuracy of independent
204 predictions increases, the model assigns a greater weight to the contextual predictions. We present
205 the hyperparameter analysis experiment for E_η in Appendix C.

206 4 EXPERIMENT

207 4.1 EXPERIMENTAL SETUP

208 **Datasets.** In this work, we use two public and one private BST data to measure the performance of
209 models. More detailed descriptions can be found in Appendix D.

- 210 • **fNIRS.** The Tufts fNIRS to Mental Workload (**Tufts fNIRS2MW** (Huang et al., 2021)) data con-
211 tains brain activity recordings and other data from adult humans performing controlled cognitive
212 workload tasks. They label each part of the experiment with one of four possible levels of n -back
213 working memory intensity. Following Huang et al. (2021), we classify 0-back and 2-back tasks.
- 214 • **Sleep.** The **SleepEDF** (Kemp et al., 2000) data contains PolySomnoGraphic sleep records for 197
215 subjects over a whole night, including EEG, EOG, chin EMG, and event markers, as well as some
216 respiration and temperature data. In our work, following Kemp et al. (2000), we use the EEG
217 Fpz-Cz channel and EOG horizontal channel.
- 218 • **SEEG.** The private SEEG data records brain signals indicative of suspected pathological tissue
219 within the brains of seizure patients. Different neurosurgeons annotate the seizure waveforms
220 within the brain signals for classification. In our work, we uniformly downsample the data to
221 250Hz and identify seizures for each single channel.

222 **Label disturbance.** We introduce a novel disturbance method to the original labels of the public
223 data to simulate scenarios where labels are inconsistent. Specifically, we first look for the boundary
224 points between different classes in a complete long time sequence. Then, we randomly determine
225 with a 0.5 probability whether each boundary point should move forward or backward. Finally, we
226 randomly select a new boundary point position from $r\%$ of the length of the class in the direction of
227 the boundary movement. In this way, we can interfere with the boundary labels and simulate label
228 inconsistency. Meanwhile, a larger value of $r\%$ indicates a higher degree of label inconsistency. In
229 this work, we conduct experiments with r values of 0, 20, and 40 for fNIRS and Sleep data.

230 **Baselines.** We compare *Con4m* with state-of-art models from various domains, including one
 231 time series classification (TSC) model with noisy labels **SREA** (Castellani et al., 2021), three image
 232 classification models with noisy labels: **SIGUA** (Han et al., 2020), **UNICON** (Karim et al., 2022)
 233 and **Sel-CL** (Li et al., 2022), one supervised TSC model **MiniRocket** (Dempster et al., 2021), one
 234 time series backbone model **TimesNet** (Wu et al., 2023), and one time series forecasting model
 235 **PatchTST** (Nie et al., 2023). See more detailed descriptions of the baselines in Appendix E.

Table 1: Overview of BST data used in this work.

Data	Sample Frequency	# of Features	# of Classes	Subjects	Groups	Cross Validation	Total Intervals	Interval Length	Window Length	Slide Length	Total Segments
fNIRS	5.2Hz	8	2	68	4	12	4,080	38.46s	4.81s	0.96s	146,880
Sleep	100Hz	2	5	154	3	6	6,000	40s	2.5s	1.25s	186,000
SEEG	250Hz	1	2	8	4	3	8,000	16s	1s	0.5s	248,000

236 **Implementation details.** We use cross-validation (Kohavi, 1995) to evaluate the model’s general-
 237 ization ability by partitioning the subjects in the data into non-overlapping subsets for training and
 238 testing. As shown in Table 1, for fNIRS and SEEG data, we divide the subjects into 4 groups and
 239 follow the 2 training-1 validation-1 testing (2-1-1) setting to conduct experiments. We divide the
 240 Sleep data into 3 groups and follow the 1-1-1 experimental setting. Therefore, we report the mean
 241 values of 12 and 6 cross-validation results for fNIRS and Sleep data respectively. Notice that for
 242 SEEG data, inconsistent labels already exist in the original data. To obtain a high-quality testing
 243 group, we select one group for accurate labeling and use a majority voting procedure to determine
 244 the boundaries. Then we leave the testing group aside and only change the validation group to report
 245 the mean value of 3 experiments. We report the full experimental results in Appendix G.

246 **Evaluation metrics.** We use Accuracy (Acc.) and Macro- F_1 (F_1) scores as our evaluation metrics
 247 due to the balanced testing set. Macro- F_1 score is the average of the F_1 scores across all classes.

248 4.2 LABEL DISTURBANCE EXPERIMENT

249 The average results over all cross-validation experiments of different methods are presented in Ta-
 250 ble 2. Overall, *Con4m* outperforms almost all baselines across all data and all disturbance ratios.

Table 2: Comparison with state-of-the-art methods in the testing Accuracy (%) and F_1 score (%) on three BST data. The **best results** are in bold and we underline the **second best results**.

Model	r%	Noisy Label Learning						Time Series Classification						Both			
		SIGUA		UNICON		Sel-CL		MiniRocket		TimesNet		PatchTST		SREA	<i>Con4m</i>		
		Acc.	F_1	Acc.	F_1	Acc.	F_1	Acc.	F_1	Acc.	F_1	Acc.	F_1	Acc.	F_1		
fNIRS	0	64.58	67.37	63.21	61.15	63.92	63.86	60.89	61.28	65.17	67.47	52.87	51.79	<u>65.18</u>	<u>70.10</u>	67.91	71.28
	20	63.45	65.24	62.33	60.45	61.85	62.45	59.74	60.41	63.48	65.39	52.42	55.38	<u>63.99</u>	<u>69.65</u>	66.78	71.27
	40	60.55	63.47	60.63	57.35	61.21	61.75	57.56	57.87	61.76	63.45	51.94	52.67	63.75	<u>69.40</u>	<u>63.50</u>	70.04
Sleep	0	54.47	54.28	62.71	62.26	<u>63.43</u>	<u>63.48</u>	62.36	62.00	59.87	59.50	58.72	58.40	49.73	48.81	67.93	68.02
	20	53.50	53.07	62.59	61.63	<u>63.19</u>	<u>63.45</u>	62.17	61.75	59.17	57.72	56.69	56.16	49.43	48.80	66.61	66.31
	40	52.16	51.32	60.65	58.34	<u>61.85</u>	<u>61.72</u>	59.19	58.38	56.68	55.73	54.21	53.05	48.22	45.72	65.34	64.31
SEEG	-	66.87	53.19	<u>69.22</u>	60.53	68.46	60.50	68.79	<u>62.39</u>	66.02	50.99	66.59	58.45	65.11	55.21	74.60	72.00

251 **Results of different methods.** For fNIRS, *Con4m* achieves competitive performance compared to
 252 SREA. SREA is particularly designed for time series data and could better identify the inconsistent
 253 time segments in a self-supervised fashion. However, on Sleep and SEEG data that require a stronger
 254 reliance on contextual information, SREA’s performance is significantly lower than *Con4m*. More-
 255 over, in the case of SEEG and Sleep data without disturbance, *Con4m* impressively improves **7.14%**
 256 and **15.41%** compared with the best baseline in F_1 score. This results demonstrate the necessity of
 257 considering contextual information when dealing with more complex independent segments.

258 **Results of different $r\%$.** Noisy label learning methods demonstrate close performance degradation
 259 as $r\%$ increases from 0% to 20%. But with a higher ratio from 20% to 40%, SIGUA, UNICON,
 260 Sel-CL and SREA show averaged 3.01%, 5.23%, 1.92% and 3.34% decrease in F_1 score across
 261 fNIRS and Sleep data, while *Con4m* shows 2.37% degradation. For TSC models, there is a consis-
 262 tent performance decline as $r\%$ rises. Non-deep learning-based MiniRocket shows a more robust

263 performance. The performance of PatchTST on fNIRS data exhibits significant instability, possibly
 264 due to its tendency to overfit inconsistent labels too quickly. The stable performance of *Con4m*
 265 indicates that our proposed training framework can effectively harmonize inconsistent labels.

266 **Results of random disturbance.** We also conduct experiments following the setting of random
 267 label disturbance, which is commonly employed in the noisy label learning works (Wei et al., 2021;
 268 Li et al., 2022; Huang et al., 2023) of the image classification domain. As shown in Figure 3(b),
 269 compared to our novel boundary disturbance, *Con4m* exhibits stronger robustness to random dis-
 270 turbance. Even with the 20% disturbance ratio, *Con4m* treats it as a form of data augmentation,
 271 resulting in improved performance. This indicates that overcoming more challenging boundary dis-
 272 turbance aligns better with the nature of time series data.

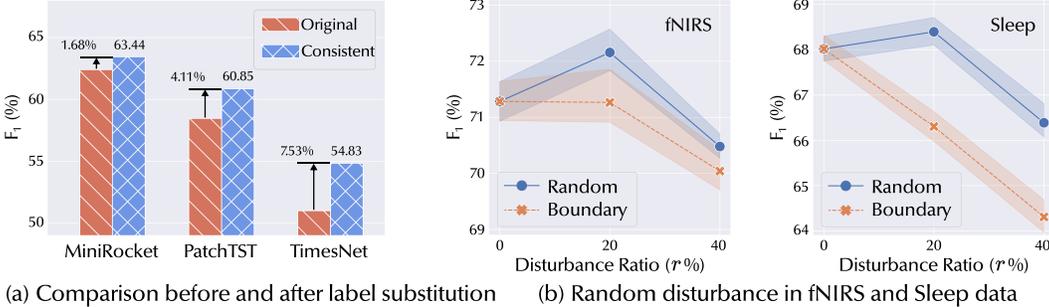


Figure 3: Comparison results of label substitution and random disturbance experiments.

273 4.3 LABEL SUBSTITUTION EXPERIMENT

274 Since blurred boundaries are inherent to SEEG data and the majority voting procedure is costly,
 275 we limit this procedure to only one high-quality testing group in the label disturbance experiment.
 276 Besides, on the SEEG data, *Con4m* modifies approximately 10% of the training labels, which is a
 277 significant proportion. Therefore, it is necessary to further evaluate the effectiveness of our label har-
 278 monization process on SEEG data. Specifically, we train the TSC baselines based on the harmonized
 279 labels generated by *Con4m* and observe to what extent the baseline results are improved. As shown
 280 in Figure 3(a), PatchTST and TimesNet, employing deep learning architectures, are more suscep-
 281 tible to label inconsistency, so they obtain more significant performance improvement (4.11% and
 282 7.53% in F_1 score). Unlike modified PatchTST that considers the classified segments in contexts,
 283 TimesNet only focuses on the independent segments, thus having a more dramatic improvement. In
 284 contrast, MiniRocket achieves only 1.68% increase. The reason may be that MiniRocket utilizes a
 285 simple random feature mapping approach without relying on specific patterns or correlations.

286 4.4 ABLATION EXPERIMENT

Table 3: Comparison with ablations in the testing Accuracy (%) and F_1 score (%) on two public data. The **best results** are in bold and we underline the second best results.

Model \ Dataset	r%	Preserve one						Remove one						<i>Con4m</i>					
		+ Con-T		+ Coh-P		+ Cur-L		- Con-T		- Coh-P		- Cur-L		- Fit		- η		Acc.	F_1
Sleep	20	65.97	65.05	65.76	65.10	65.31	64.76	65.73	<u>65.53</u>	65.84	65.07	65.85	65.43	<u>66.06</u>	65.28	62.02	59.97	66.61	66.31
	40	63.94	62.67	64.42	62.76	63.69	62.23	64.44	63.05	64.23	63.03	<u>64.89</u>	63.07	64.69	<u>63.22</u>	61.93	57.98	65.34	64.31
SEEG	-	71.68	67.85	71.69	69.04	71.32	67.22	73.85	70.59	72.41	68.26	<u>74.17</u>	<u>71.18</u>	73.47	70.63	70.70	66.04	74.60	72.00

287 We introduce two types of model variations. **(1) Preserve one module.** We preserve only the
 288 Con-Transformer (Con-T), Coherent Prediction (Coh-P), or Curriculum Learning (Cur-L) module
 289 separately. **(2) Remove one component.** In addition to removing the above three modules, we also
 290 remove the function fitting component (-Fit) and η ($E_\eta = 0$) to verify the necessity of prediction
 291 behavior constraint and progressively updating labels.

292 As shown in Table 3, when keeping one module, +Coh-P achieves the best performance with aver-
 293 aged 2.78% decrease in F_1 score, indicating that introducing the contextual class information are
 294 most effective for BST data. The utility of each module varies across datasets. For example, for
 295 Sleep data, the Con-T contributes more to performance improvement compared to the Cur-L mod-
 296 ule, while the opposite phenomenon is observed for SEEG data. As for removing one component,
 297 even when we only remove the Tanh function fitting, the F_1 score of *Con4m* significantly decreases
 298 1.72% on average. On the Sleep-20% and SEEG data, the drop caused by -Fit is more significant
 299 than that caused by some other modules. Moreover, the model variation $-\eta$ achieves the worst results
 300 (9.23% F_1 drop), aligning with our motivation. Specifically, during early training stages, the model
 301 tends to learn the consistent parts of the original labels. Premature use of unreliable predicted labels
 302 as subsequent training supervision signals leads to model poisoning and error accumulation.

303 4.5 CASE STUDY

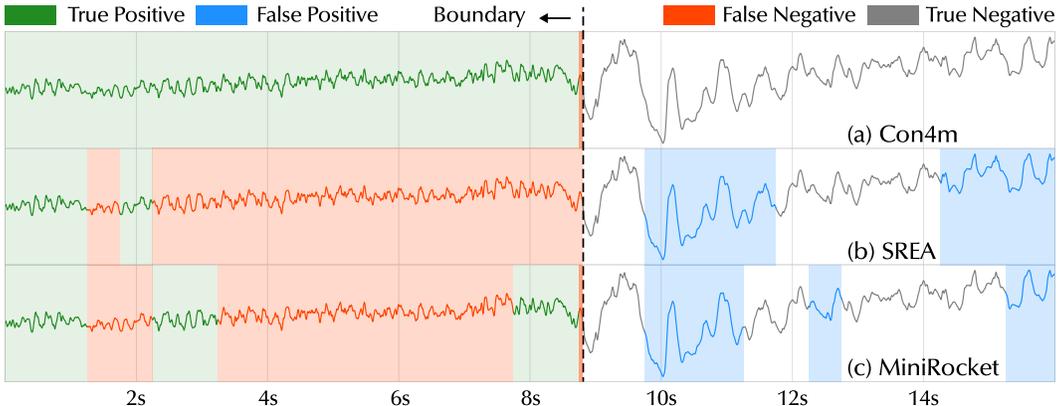


Figure 4: Case study for a continuous time interval in SEEG testing set.

304 We present a case study to provide a specific example that illustrates how *Con4m* works for BST
 305 data in Figure 4. We show a comparative visualization result of *Con4m*, SREA and MiniRocket for
 306 the predictions in a continuous time interval in SEEG testing set. In SEEG data, we assign the label
 307 of normal segments as 0 and that of seizures as 1. As the figure shows, *Con4m* demonstrates a more
 308 coherent narrative by constraining the prediction behavior to align with the contextual information of
 309 data. In contrast, MiniRocket and SREA exhibit noticeably interrupted and inconsistent predictions.
 310 What is even more impressive is that the model accurately identifies consistent boundaries within the
 311 time intervals spanning across two different states. This verifies that the harmonized labels capture
 312 the boundaries between distinct classes more precisely. Refer to Appendix H for more cases.

313 5 CONCLUSION AND DISCUSSION

314 In this work, we introduce the conception of Blurred-Segmented Time Series (BST) data and pose
 315 its unique challenges which have been overlooked by mainstream time series classification (TSC)
 316 models. Through theoretical analysis, we have obtained the conclusion that valuable contextual in-
 317 formation enhances the predictive ability of the model. By introducing a novel method, *Con4m*,
 318 we incorporate effective contextual information at both the data and class levels to enhance model’s
 319 predictive ability. Extensive experiments not only validate the superior performance achieved by
 320 *Con4m* through the integration of valuable contextual information, but also highlight the effective-
 321 ness and necessity of the proposed consistent label training framework for modeling BST data. Our
 322 approach still has some limitations. We have solely focused on analyzing and designing end-to-end
 323 supervised models. Further exploration to self-supervised methods would be challenging yet in-
 324 triguing. When faced with more diverse label behaviors, the function fitting module needs to engage
 325 in more selection and design of basis functions. Nevertheless, our work brings new insights to the
 326 classification-based fields. In particular, for the TSC domain, we re-emphasize the importance of the
 327 inherent temporal dependence of time segments, shedding light on the era of personalized services.

328 REFERENCES

- 329 Dana Angluin and Philip Laird. Learning from noisy examples. *Mach. Learn.*, 2(4):343—370,
330 1988.
- 331 Samaneh Azadi, Jiashi Feng, Stefanie Jegelka, and Trevor Darrell. Auxiliary image regularization
332 for deep cnns with noisy labels. In *International Conference on Learning Representations (ICLR)*,
333 2016.
- 334 Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul
335 Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018. *arXiv*
336 *preprint arXiv:1811.00075*, 2018.
- 337 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
338 learning to align and translate. In *International Conference on Learning Representations (ICLR)*,
339 2015.
- 340 Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds.
341 *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- 342 Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In
343 *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, pp. 41–
344 –48, 2009.
- 345 Andrea Castellani, Sebastian Schmitt, and Barbara Hammer. Estimating the electrical power output
346 of industrial devices with end-to-end time-series classification in the presence of label noise. In
347 *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML*
348 *PKDD)*, pp. 469–484, 2021.
- 349 Junru Chen, Yang Yang, Tao Yu, Yingying Fan, Xiaolong Mo, and Carl Yang. Brainnet: Epileptic
350 wave detection from seeg with hierarchical graph diffusion learning. In *Proceedings of the 28th*
351 *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 2741–2751,
352 2022.
- 353 Ranak Roy Chowdhury, Xiyuan Zhang, Jingbo Shang, Rajesh K. Gupta, and Dezhi Hong. Tarnet:
354 Task-aware reconstruction for time-series transformer. In *Proceedings of the 28th ACM SIGKDD*
355 *Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 212—220, 2022.
- 356 Hoang Anh Dau, Anthony Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh
357 Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn Keogh. The ucr time series archive.
358 *IEEE/CAA Journal of Automatica Sinica*, 6(6):1293–1305, 2019.
- 359 Angus Dempster, François Petitjean, and Geoffrey I. Webb. Rocket: Exceptionally fast and accurate
360 time series classification using random convolutional kernels. *Data Min. Knowl. Discov.*, 34(5):
361 1454—1495, 2020.
- 362 Angus Dempster, Daniel F. Schmidt, and Geoffrey I. Webb. Minirocket: A very fast (almost) de-
363 terministic transform for time series classification. In *Proceedings of the 27th ACM SIGKDD*
364 *Conference on Knowledge Discovery & Data Mining (KDD)*, pp. 248—257, 2021.
- 365 Don Dennis, Durmus Alp Emre Acar, Vikram Mandikal, Vinu Sankar Sadasivan, Venkatesh
366 Saligrama, Harsha Vardhan Simhadri, and Prateek Jain. Shallow rnn: Accurate time-series classi-
367 fication on resource constrained devices. In *Advances in Neural Information Processing Systems*
368 *(NeurIPS)*, 2019.
- 369 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
370 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference*
371 *of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pp.
372 4171–4186, 2019.
- 373 Arsalan Dezhkam, Mohammad Taghi Manzuri, Ahmad Aghapour, Afshin Karimi, Ali Rabiee, and
374 Shervin Manzuri Shalmani. A bayesian-based classification framework for financial time series
375 trend prediction. *J. Supercomput.*, 79(4):4622—4659, 2022.

- 376 Navid Mohammadi Foumani, Lynn Miller, Chang Wei Tan, Geoffrey I Webb, Germain Forestier,
377 and Mahsa Salehi. Deep learning for time series classification and extrinsic regression: A current
378 survey. *arXiv preprint arXiv:2302.02515*, 2023.
- 379 Vivien Sainte Fare Garnot and Loic Landrieu. Lightweight temporal self-attention for classifying
380 satellite images time series. In *Advanced Analytics and Learning on Temporal Data: 5th ECML*
381 *PKDD Workshop (AALTD)*, pp. 171—181, 2020.
- 382 Tieliang Gong, Qian Zhao, Deyu Meng, and Zongben Xu. Why curriculum learning & self-paced
383 learning work in big/noisy data: A theoretical perspective. *Big Data and Information Analytics*,
384 1(1):111–127, 2016.
- 385 Alex Graves, Marc G. Bellemare, Jacob Menick, Rémi Munos, and Koray Kavukcuoglu. Automated
386 curriculum learning for neural networks. In *Proceedings of the 34th International Conference on*
387 *Machine Learning (ICML)*, pp. 1311–1320, 2017.
- 388 Bo Han, Jiangchao Yao, Gang Niu, Mingyuan Zhou, Ivor Tsang, Ya Zhang, and Masashi Sugiyama.
389 Masking: A new perspective of noisy supervision. In *Advances in Neural Information Processing*
390 *Systems (NeurIPS)*, 2018a.
- 391 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
392 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In
393 *Advances in Neural Information Processing Systems (NeurIPS)*, 2018b.
- 394 Bo Han, Gang Niu, Xingrui Yu, Quanming Yao, Miao Xu, Ivor Tsang, and Masashi Sugiyama.
395 SIGUA: Forgetting may make learning with noisy labels more robust. In *Proceedings of the 37th*
396 *International Conference on Machine Learning (ICML)*, pp. 4006–4016, 2020.
- 397 Bo Han, Quanming Yao, Tongliang Liu, Gang Niu, Ivor W. Tsang, James T. Kwok, and Masashi
398 Sugiyama. A survey of label-noise representation learning: Past, present and future. *arXiv*
399 *preprint arXiv:2011.04406*, 2021.
- 400 Zhe Huang, Liang Wang, Giles Blaney, Christopher Slaughter, Devon McKeon, Ziyu Zhou, Robert
401 Jacob, Robert Jacob, and Michael Hughes. The tufts fnirs mental workload dataset benchmark for
402 brain-computer interfaces that generalize. In *Proceedings of the Neural Information Processing*
403 *Systems (NeurIPS) Track on Datasets and Benchmarks*, 2021.
- 404 Zhizhong Huang, Junping Zhang, and Hongming Shan. Twin contrastive learning with noisy la-
405 bels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
406 *(CVPR)*, pp. 11661–11670, 2023.
- 407 Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain
408 Muller. Deep learning for time series classification: A review. *Data Min. Knowl. Discov.*, 33(4):
409 917—963, 2019.
- 410 Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning
411 data-driven curriculum for very deep neural networks on corrupted labels. *arXiv preprint*
412 *arXiv:1712.05055*, 2018.
- 413 Nazmul Karim, Mamshad Nayeem Rizve, Nazanin Rahnavard, Ajmal Mian, and Mubarak Shah.
414 Unicon: Combating label noise through uniform selection and contrastive learning. In *Proceed-*
415 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
416 9676–9686, 2022.
- 417 B. Kemp, A.H. Zwinderman, B. Tuk, H.A.C. Kamphuisen, and J.J.L. Obery. Analysis of a sleep-
418 dependent neuronal feedback loop: The slow-wave microcontinuity of the eeg. *IEEE Transactions*
419 *on Biomedical Engineering*, 47(9):1185–1194, 2000.
- 420 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
421 *Conference on Learning Representations (ICLR)*, 2015.
- 422 Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection.
423 In *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*, pp.
424 1137—1143, 1995.

- 425 M. Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In
426 *Advances in Neural Information Processing Systems (NeurIPS)*, 2010.
- 427 Neil D. Lawrence and Bernhard Schölkopf. Estimating a kernel fisher discriminant in the presence
428 of label noise. In *Proceedings of the Eighteenth International Conference on Machine Learning*
429 (*ICML*), pp. 306—313, 2001.
- 430 Shikun Li, Xiaobo Xia, Shiming Ge, and Tongliang Liu. Selective-supervised contrastive learning
431 with noisy labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
432 *Recognition (CVPR)*, pp. 316–325, 2022.
- 433 Jason Lines, Sarah Taylor, and Anthony Bagnall. Time series classification with hive-cote: The
434 hierarchical vote collective of transformation-based ensembles. *ACM Trans. Knowl. Discov. Data*,
435 12(5), 2018.
- 436 Bing Liu and Huanhuan Cheng. Financial time series classification method based on low-frequency
437 approximate representation. *Engineering Reports*, pp. e12739, 2023.
- 438 Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. Robust training under label noise by over-
439 parameterization. In *Proceedings of the 39th International Conference on Machine Learning*
440 (*ICML*), pp. 14153–14172, 2022.
- 441 Tabet Matiisen, Avital Oliver, Taco Cohen, and John Schulman. Teacher–student curriculum learn-
442 ing. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 31(9):3732–3740,
443 2020.
- 444 Miller McPherson, Lynn Smith-Lovin, and James M Cook. Birds of a feather: Homophily in social
445 networks. *Annual review of sociology*, 27(1):415–444, 2001.
- 446 Matthew Middlehurst, Patrick Schäfer, and Anthony Bagnall. Bake off redux: A review
447 and experimental evaluation of recent time series classification algorithms. *arXiv preprint*
448 *arXiv:2304.13029*, 2023.
- 449 Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64
450 words: Long-term forecasting with transformers. In *The Eleventh International Conference on*
451 *Learning Representations (ICLR)*, 2023.
- 452 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
453 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
454 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
455 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep
456 learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 457 Giorgio Patrini, Alessandro Rozza, Aditya Menon, Richard Nock, and Lizhen Qu. Making deep neu-
458 ral networks robust to label noise: a loss correction approach. *arXiv preprint arXiv:1609.03683*,
459 2017.
- 460 Alireza Rafiei, Rasoul Zahedifar, Chiranjibi Sitaula, and Faezeh Marzbanrad. Automated detection
461 of major depressive disorder with eeg signals: A time series classification using deep learning.
462 *IEEE Access*, 10:73804–73817, 2022.
- 463 Deepta Rajan and Jayaraman J. Thiagarajan. A generative modeling approach to limited channel
464 ecg classification. In *2018 40th Annual International Conference of the IEEE Engineering in*
465 *Medicine and Biology Society (EMBC)*, pp. 2571–2574, 2018.
- 466 Zezhi Shao, Zhao Zhang, Fei Wang, and Yongjun Xu. Pre-training enhanced spatial-temporal graph
467 neural network for multivariate time series forecasting. In *Proceedings of the 28th ACM SIGKDD*
468 *Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 1567—1577, 2022.
- 469 Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training
470 convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2015.

- 471 Sijie Tian, Yaoyu Zhang, Yuchun Feng, Nour Elsagan, Yoon Ko, M. Hamed Mozaffari, Dexten D.Z.
472 Xi, and Chi-Guhn Lee. Time series classification, augmentation and artificial-intelligence-enabled
473 software for emergency response in freight transportation fires. *Expert Systems with Applications*,
474 233:120914, 2023.
- 475 Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. *IEEE Transactions*
476 *on Pattern Analysis and Machine Intelligence (TPAMI)*, 44(9):4555–4576, 2021.
- 477 Hongxin Wei, Lue Tao, RENCHUNZI XIE, and Bo An. Open-set label noise can improve robustness
478 against inherent label noise. In *Advances in Neural Information Processing Systems (NeurIPS)*,
479 pp. 7978–7992, 2021.
- 480 Daphna Weinshall, Gad Cohen, and Dan Amir. Curriculum learning by transfer learning: Theory
481 and experiments with deep networks. In *Proceedings of the 35th International Conference on*
482 *Machine Learning (ICML)*, pp. 5238–5246, 2018.
- 483 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
484 Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International*
485 *Conference on Learning Representations (ICLR)*, 2023.
- 486 Kai Wu, Kaixin Yuan, Yingzhi Teng, Jing Liu, and Licheng Jiao. Broad fuzzy cognitive map systems
487 for time series classification. *Appl. Soft Comput.*, 128(C):109458, 2022.
- 488 Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series
489 anomaly detection with association discrepancy. In *International Conference on Learning Repre-*
490 *sentations (ICLR)*, 2022.
- 491 Chao-Han Huck Yang, Yun-Yun Tsai, and Pin-Yu Chen. Voice2series: Reprogramming acoustic
492 models for time series classification. In *Proceedings of the 38th International Conference on*
493 *Machine Learning (ICML)*, pp. 11808–11819, 2021.
- 494 Chong You, Zhihui Zhu, Qing Qu, and Yi Ma. Robust recovery via implicit bias of discrepant
495 learning rates for double over-parameterization. *arXiv preprint arXiv:2006.08857*, 2020.
- 496 Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W. Tsang, and Masashi Sugiyama. How does
497 disagreement help generalization against label corruption? *arXiv preprint arXiv:1901.04215*,
498 2019.
- 499 Yuan Yuan, Lei Lin, Qingshan Liu, Renlong Hang, and Zeng-Guang Zhou. Sits-former: A pre-
500 trained spatio-spectral-temporal representation model for sentinel-2 time series classification. *In-*
501 *ternational Journal of Applied Earth Observation and Geoinformation*, 106:102651, 2022.
- 502 Bowen Zhao, Huanlai Xing, Xinhan Wang, Fuhong Song, and Zhiwen Xiao. Rethinking attention
503 mechanism in time series classification. *Inf. Sci.*, 627(C):97–114, 2023.
- 504 Evgenii Zheltonozhskii, Chaim Baskin, Avi Mendelson, Alex M. Bronstein, and Or Litany. Con-
505 trast to divide: Self-supervised pre-training for learning with noisy labels. In *Proceedings of*
506 *the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 1657–1667,
507 2022.
- 508 Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond
509 homophily in graph neural networks: Current limitations and effective designs. In *Advances in*
510 *Neural Information Processing Systems (NeurIPS)*, pp. 7793–7804, 2020.

511 A DETAILS OF RELATED WORKS

512 **Time series classification (TSC).** TSC has become a popular field in various applications with
 513 the exponential growth of available time series data in recent years. In response, researchers have
 514 proposed numerous algorithms (Ismail Fawaz et al., 2019). High accuracy in TSC is achieved by
 515 classical algorithms such as Rocket and its variants (Dempster et al., 2020; 2021), which use random
 516 convolution kernels with relatively low computational cost, as well as ensemble methods like HIVE-
 517 COTE (Lines et al., 2018), which assign weights to individual classifiers.

518 Moreover, the flourishing non-linear modeling capacity of deep models has led to an increasing
 519 prevalence of TSC algorithms based on deep learning. Various techniques are utilized in TSC:
 520 RNN-based methods (Rajan & Thiagarajan, 2018; Dennis et al., 2019) capture temporal changes
 521 through state transitions; MLP-based methods (Garnot & Landrieu, 2020; Wu et al., 2022) encode
 522 temporal dependencies into parameters of the MLP layer; and the latest method TimesNet (Wu et al.,
 523 2023) converts one-dimensional time series into a two-dimensional space, achieving state-of-the-art
 524 performance on five mainstream tasks. Furthermore, Transformer-based models (Yang et al., 2021;
 525 Chowdhury et al., 2022) with attention mechanism have been widely used.

526 The foundation of our work lies in these researches, including the selection of the backbone and
 527 experimental setup. However, mainstream TSC models (Middlehurst et al., 2023; Foumani et al.,
 528 2023) are often designed for publicly available datasets (Bagnall et al., 2018; Dau et al., 2019) based
 529 on the *i.i.d.* samples, disregarding the inherent contextual dependencies between classified samples
 530 in BST data. Although some time series models (Shao et al., 2022; Nie et al., 2023) use patch-by-
 531 patch technique to include contextual information, they are partially context-aware since they only
 532 model the data dependencies between time points, ignoring the class dependencies of segments.

533 **Noisy label learning (NLL).** NLL is an important and challenging research topic in machine learn-
 534 ing, as real-world data often rely on manual annotations prone to errors. Early works focus on
 535 statistical learning (Angluin & Laird, 1988; Lawrence & Schölkopf, 2001; Bartlett et al., 2006). Re-
 536 searches including Sukhbaatar et al. (2015) launch the era of noise-labeled representation learning.

537 The label noise transition matrix, which represents the transition probability from clean labels to
 538 noisy labels (Han et al., 2021), is an essential tool. Common techniques for loss correction include
 539 forward and backward correction (Patrini et al., 2017), while masking invalid class transitions with
 540 prior knowledge is also an important method (Han et al., 2018a). Adding an explicit or implicit
 541 regularization term in objective functions can reduce the model’s sensitivity to noise, whereas re-
 542 weighting mislabeled data can reduce its impact on the objective (Azadi et al., 2016; You et al., 2020;
 543 Liu et al., 2022). Other methods involve training on small-loss instances and utilizing memorization
 544 effects. MentorNet (Jiang et al., 2018) pretrains a secondary network to choose clean instances
 545 for primary network training. Co-teaching (Han et al., 2018b) and Co-teaching+ (Yu et al., 2019),
 546 as sample selection methods, introduce two neural networks with differing learning capabilities to
 547 train simultaneously, which filter noise labels mutually. The utilization of contrastive learning has
 548 emerged as a promising approach for enhancing the robustness in the context of classification tasks
 549 of label correction methods (Li et al., 2022; Zheltonozhskii et al., 2022; Huang et al., 2023).

550 These works primarily focus on handling noisy labels. And ensuring overall label consistency by
 551 modifying certain labels is crucial for BST data. To the best of our knowledge, the only noisy label
 552 learning work in the time series field is SREA (Castellani et al., 2021), which trains a classifier and
 553 an autoencoder with a shared embedding representation, progressively self-relabeling mislabeled
 554 data samples in a self-supervised manner. However, SREA does not take into account the contextual
 555 dependencies of BST data, limiting its performance.

556 **Curriculum learning (CL).** Bengio et al. (2009) propose CL, which imitates human learning by
 557 starting with simple samples and progressing to complicated ones. Based on this notion, CL can
 558 denoise noisy data since learners are encouraged to train on easier data and spend less time on noisy
 559 samples (Gong et al., 2016; Wang et al., 2021). Current mainstream approaches include Self-paced
 560 Learning (Kumar et al., 2010), where students schedule their learning, Transfer Teacher (Weinshall
 561 et al., 2018), based on a predefined training scheduler; and RL Teacher (Graves et al., 2017; Matiisen
 562 et al., 2020), which incorporates student feedback into the framework. The utilization of CL proves
 563 to be particularly advantageous in situations involving changes in the training labels. Hence, this
 564 technique is utilized to enhance the modeling of BST data in a more stable manner.

565 B IMPLEMENTATION DETAILS OF PREDICTION BEHAVIOR CONSTRAINT

566 To fit the hyperbolic tangent function (Tanh), we use the mean squared error (MSE) loss function. In
 567 practice, we use the Adam optimizer with a learning rate of 0.1 to optimize the trainable parameters.
 568 The maximum number of iterations is set to 100, and the tolerance value for stopping the fitting
 569 process based on loss change is set to $1e-6$. Sequences belonging to one minibatch are parallelized
 570 to fit their respective Tanh functions. To adapt to the value range of the standard Tanh function, we
 571 rescale the sequential predictions to $[-1, 1]$ before fitting.

572 However, it can be difficult to achieve a good fit when fitting with the Tanh function. Specifically,
 573 random initialization may fail to fit the sequential values properly when a long time series undergoes
 574 a state transition near the boundary. For example, as Figure 5(a) shows, we fit a sequence in which
 575 only the last value is 1. We set all default initial parameters as 1 and fit it. It can be observed that
 576 the fitting function cannot properly fit the trend and will mislabel the last point.

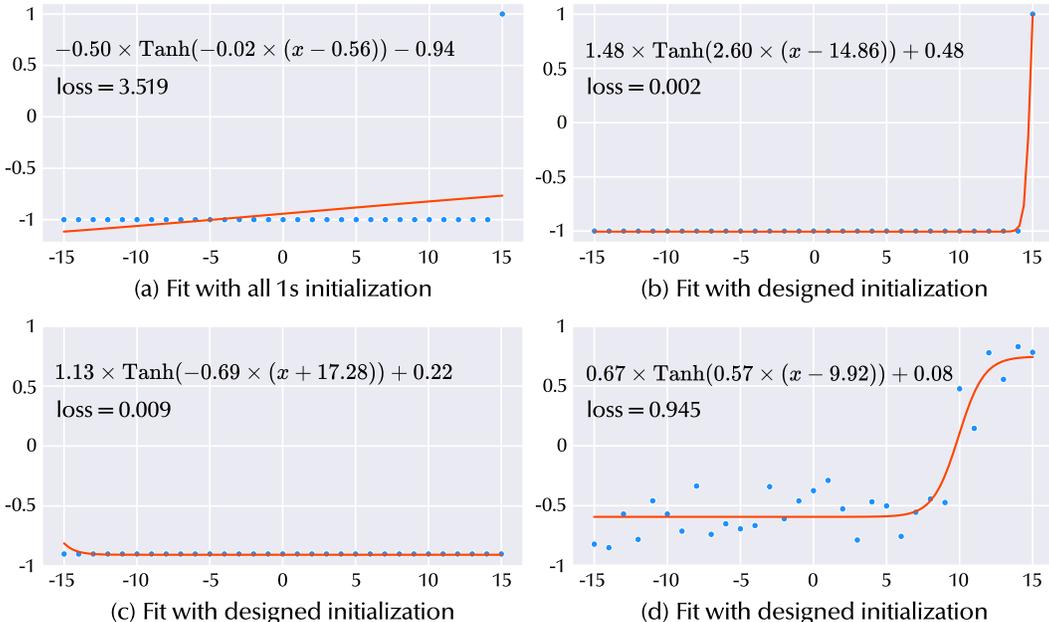


Figure 5: Cases for Tanh fitting.

577 Appropriate parameter initialization is needed to avoid excessive bias. After careful observation,
 578 we find that parameter k controls the slope at the transition part of Tanh, and parameter b controls
 579 the abscissa at the transition point. In the process, all fitting values are assigned with uniform
 580 abscissa values. Therefore, we calculate the maximum difference between adjacent values and the
 581 corresponding position in the entire sequence. And these two values are assigned to parameters k
 582 and b , respectively. This allows us to obtain suitable initial parameters and avoid getting trapped in
 583 local optima or saddle points during function fitting. Formally, given the L -length input sequence \tilde{p} ,
 584 we initialize parameters k and b as follows:

$$di = [\tilde{p}_{i+1} - \tilde{p}_i]_{i \in \{1, \dots, L-1\}}, \quad (19)$$

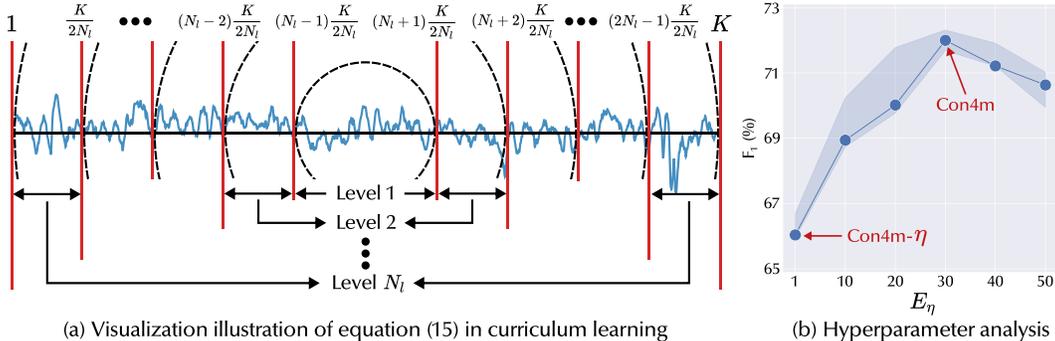
$$k, b = \max(\text{Abs}(di)), \arg \max(\text{Abs}(di)), \quad (20)$$

$$k = k \times \text{Sign}(di[b]), \quad (21)$$

$$b = -(b - \lfloor L/2 \rfloor + 0.5), \quad (22)$$

585 where $\text{Abs}(\cdot)$ and $\text{Sign}(\cdot)$ denote the absolute value function and sign function respectively. di is the
 586 difference vector. After proper initialization, as Figure 5(b) shows, we can obtain more accurate fit-
 587 ting results to reduce the probability of mislabeling. We also show some other cases (Figure 5(c)(d))
 588 for the fitting results to verify the effectiveness of the fitting process we propose.

589 C HYPERPARAMETER ANALYSIS

Figure 6: Visualization of data division in curriculum learning and hyperparameter analysis of E_η .

590 The dynamic weighting factor η is introduced to progressively update the labels, preventing the
 591 model from overly relying on its own predicted labels too early. To validate the utility of η and
 592 determine an appropriate linear growth epoch E_η , we conduct the hyperparameter search experiment
 593 on SEEG data. As shown in Figure 6(b), with smaller E_η (corresponding to a higher growth rate),
 594 there is a significant improvement in model performance. This aligns with our motivation that during
 595 the early stage of model training, the primary objective is to better fit the original labels. At this
 596 stage, the model’s own predictions are unreliable. If the predicted results are used as training labels
 597 too early in subsequent epochs, the model would be adversely affected by its own unreliability.
 598 On the other hand, excessively large E_η leads to a slower rate of label updates, making it more
 599 challenging for the model to timely harmonize inconsistent labels. Nonetheless, considering the
 600 impact of variance, the model exhibits robustness to slightly larger E_η . In this work, we uniformly
 601 use $E_\eta = 30$ as the default value.

602 D DETAILS OF DATASETS

603 **fNIRS.** All signals are sampled at a frequency of 5.2Hz. At each time step, they record 8 real-valued
 604 measurements, with each measurement corresponding to 2 concentration changes (oxyhemoglobin
 605 and deoxyhemoglobin), 2 types of optical data (intensity and phase), and 2 spatial positions on the
 606 forehead. Each measurement unit is a micromolar concentration change per liter of tissue (for oxy-
 607 /deoxyhemoglobin). They label each part of the active experiment with one of four possible levels
 608 of n -back working memory intensity (0-back, 1-back, 2-back, or 3-back). More specifically, in an
 609 n -back task, the subject receives 40 numbers in sequence. If a number matches the number n steps
 610 back, the subject is required to respond accordingly. There are 16 rounds of tasks, with a 20-second
 611 break between each task. Following Huang et al. (2021), we only apply classification tasks for 0-
 612 back and 2-back tasks in our work. Therefore, we only extract sequences for 0-back and 2-back
 613 tasks and concatenate them in chronological order.

614 **Sleep.** The Sleep-EDF database records PolySomnoGraphic sleep data from 197 subjects, including
 615 EEG, EOG, chin EMG, and event markers. Some data also includes respiration and temperature-
 616 related signals. The database contains two studies: the Sleep Cassette study and the Sleep Telemetry
 617 study. The former records approximately 40 hours of sleep from two consecutive nights, while the
 618 latter records around 18 hours of sleep. Well-trained technicians manually score the corresponding
 619 sleep graphs according to the Rechtschaffen and Kales manual. The data is labeled in intervals of 30
 620 seconds, with each interval being marked as one of the eight possible stages: W, R, 1, 2, 3, 4, M, or ?.
 621 In our work, we utilize only the data from the Sleep Cassette study, and retain only the signals from
 622 the EEG Fpz-Cz channel and EOG horizontal channel. The EEG and EOG signals were sampled at
 623 a frequency of 100Hz. Following Kemp et al. (2000), we remove the labels for stages ? and M from
 624 the data, and merge stages 3 and 4, resulting in a 5-classification task.

625 **SEEG.** The private SEEG data records brain signals indicative of suspected pathological tissue
 626 within the brains of seizure patients. They are anonymously collected from a top hospital we coop-

627 erate with. For a patient suffering from epilepsy, 4 to 11 invasive electrodes with 52 to 153 channels
 628 are used for recording signals. In total, we have collected 847 hours of SEEG signals with a high
 629 frequency (1,000Hz or 2,000Hz) and a total capacity of 1.2TB. Professional neurosurgeons help
 630 us label the seizure segments for each channel. Before sampling for the database, we remove the
 631 bad channels marked by neurosurgeons. Then we uniformly downsample the data to 250Hz and
 632 use a low-pass filter to process the data with a cutoff frequency of 30Hz. Finally, we normalize and
 633 sample the intervals for each channel respectively.

634 E IMPLEMENTATION DETAILS OF BASELINES

- 635 • **SREA** (Castellani et al., 2021): This time series classification model with noisy labels jointly
 636 trains a classifier and an autoencoder with shared embedding representations. It gradually corrects
 637 the mislabelled data samples during training in a self-supervised fashion. We use the default model
 638 architecture from the source code provided by the author (<https://github.com/Castel44/SREA>).
- 639 • **SIGUA** (Han et al., 2020): This model adopts gradient descent on good data as usual, and learning-
 640 rate-reduced gradient ascent on bad data, thereby trying to reduce the effect of noisy labels. We
 641 modify the network for time series data based on the open source code provided by SREA, using
 642 the code from the author (<https://github.com/bhanML/SIGUA>).
- 643 • **UNICON** (Karim et al., 2022): UNICON introduces a Jensen-Shannon divergence-based uniform
 644 selection mechanism and uses contrastive learning to further combat the memorization of noisy
 645 labels. We modify the model for time series data according to the code provided by the author
 646 (<https://github.com/nazmul-karim170/UNICON-Noisy-Label>).
- 647 • **Sel-CL** (Li et al., 2022): Selective-Supervised Contrastive Learning (Sel-CL) is a latest baseline
 648 model in the field of computer vision. It selects confident pairs out of noisy ones for supervised
 649 contrastive learning (Sup-CL) without knowing noise rates. We modify the code for time series
 650 data, based on the source code provided by the author (<https://github.com/ShikunLi/Sel-CL>).
- 651 • **MiniRocket** (Dempster et al., 2021): Rocket (Dempster et al., 2020) achieves state-of-the-art ac-
 652 curacy for time series classification by transforming input time series using random convolutional
 653 kernels, and using the transformed features to train a linear classifier. MiniRocket is a variant of
 654 Rocket that improves processing time, while offering essentially the same accuracy. We use the
 655 code interface from the sktime package (<https://github.com/sktime/sktime>).
- 656 • **PatchTST** (Nie et al., 2023): This is a self-supervised representation learning framework for mul-
 657 tivariate time series by segmenting time series into subseries level patches, which are served as
 658 input tokens to Transformer with channel-independence. We modify the code to achieve classi-
 659 fication for each patch, based on the source code from the Time Series Library (TSlib) package
 660 (<https://github.com/thuml/Time-Series-Library>).
- 661 • **TimesNet** (Wu et al., 2023): This model focuses on temporal variation modeling. With Times-
 662 Block, it can discover the multi-periodicity adaptively and extract the complex temporal variations
 663 from transformed 2D tensors by a parameter-efficient inception block. We use the open source
 664 code from the TSlib package (<https://github.com/thuml/Time-Series-Library>).

665 F IMPLEMENTATION DETAILS OF *Con4m*

666 The non-linear encoder g_{enc} used in *Con4m* is composed of three 1-D convolution layers. The num-
 667 ber of kernels vary across different data and you can find corresponding parameters in the default
 668 config file of our source code. We construct the Con-Transformer based on the public codes imple-
 669 mented by HuggingFace¹. We set $d = 128$ and the dimension of intermediate representations in
 670 FFN module as 256 for all experiments. The number of heads and dropout rate are set as 8 and 0.1
 671 respectively. Since we observe that one-layer Con-Attention can fit the data well, we do not stack
 672 more layers to avoid overfitting. Note that *Con4m* consists of two Con-Transformers, we indeed use
 673 two Con-Attention layers. The model is optimized using Adam optimizer (Kingma & Ba, 2015)
 674 with a learning rate of $1e - 3$ and weight decay of $1e - 4$, and the batch size is set as 64. Also,
 675 we build our model using PyTorch 2.0.0 (Paszke et al., 2019) with CUDA 11.8. And the model is

¹https://github.com/huggingface/transformers/blob/v4.25.1/src/transformers/models/bert/modeling_bert.py

676 trained on a workstation (Ubuntu system 20.04.5) with 2 CPUs (AMD EPYC 7H12 64-Core Proces-
 677 sor) and 8 GPUs (NVIDIA GeForce RTX 3090). You can find more technical details in our source
 678 code, which has been attached in the supplementary materials.

679 G FULL RESULTS

680 The full results of the label disturbance experiment are listed in Table 4, 5 and 6. For fNIRS, we
 681 first divide the data into 4 groups by subjects and follow the 2 training-1 validation-1 testing (2-1-1)
 682 setting to conduct cross-validation experiments. Therefore, there are $C_4^2 \times C_2^1 = 12$ experiments in
 683 total. Similarly, we divide the Sleep data into 3 groups and follow the 1-1-1 experimental setting.
 684 Therefore, we carry out $C_3^1 \times C_2^1 = 6$ experiments. For SEEG data, we follow the same setting as
 685 fNIRS. Notice that we only select one group for accurate labeling to obtain a high-quality testing
 686 group, so we only have $C_3^2 = 3$ experiments. All the experimental results are listed in lexicograph-
 687 ical order according to the group name composition. We also report the mean value and standard
 688 derivation of experiments for each data.

Table 4: Full results of the label disturbance experiment on fNIRS data. The **best results** are in bold and we underline the second best results.

		Noisy Label Learning						Time Series Classification						Both			
		SIGUA		UNICON		Sel-CL		MiniRocket		TimesNet		PatchTST		SREA		Con4m	
r%	Exp	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁
0	1	62.01	64.75	62.18	63.85	63.06	63.95	60.89	61.37	61.14	60.73	51.61	51.07	65.06	69.13	<u>64.61</u>	<u>68.55</u>
	2	63.07	67.55	62.81	57.17	64.45	65.49	60.15	62.45	<u>65.11</u>	68.25	53.04	48.21	64.56	<u>70.29</u>	66.07	71.64
	3	62.93	65.56	60.16	61.71	63.67	61.44	60.81	60.96	<u>63.99</u>	66.38	51.70	54.23	<u>65.66</u>	<u>67.14</u>	66.20	70.51
	4	65.22	68.83	64.46	60.74	62.71	63.23	61.19	62.24	<u>67.42</u>	69.73	52.79	55.07	65.46	<u>70.23</u>	67.85	72.65
	5	63.28	67.96	61.70	54.87	61.46	63.21	60.46	61.35	63.07	66.46	52.99	57.11	<u>64.48</u>	70.13	66.54	<u>68.55</u>
	6	<u>65.95</u>	<u>70.12</u>	64.31	58.96	64.59	64.66	60.83	61.79	64.35	69.58	53.72	57.22	64.91	69.66	68.97	72.99
	7	61.14	64.03	60.74	62.54	61.85	60.13	60.38	60.11	61.72	62.64	52.44	50.46	<u>64.89</u>	<u>70.55</u>	68.34	70.63
	8	67.40	69.15	64.04	63.05	67.46	65.31	64.08	62.90	<u>68.78</u>	69.57	53.93	49.75	65.65	<u>70.90</u>	70.52	73.36
	9	<u>65.76</u>	68.41	63.12	66.77	61.52	63.45	58.60	58.78	64.31	67.23	51.97	48.86	64.98	<u>70.29</u>	69.37	72.60
	10	<u>68.10</u>	68.24	66.45	59.40	66.41	65.47	61.61	60.75	69.16	<u>71.17</u>	54.36	55.63	66.12	71.59	67.53	69.38
	11	64.53	65.95	63.24	57.84	<u>64.84</u>	65.24	59.12	60.71	64.18	66.64	51.65	48.51	63.71	<u>69.17</u>	65.75	69.42
	12	65.62	67.89	65.36	66.88	<u>65.02</u>	64.70	62.50	61.93	<u>68.75</u>	71.30	54.26	45.40	66.72	<u>72.17</u>	73.12	75.14
	Avg	64.58	67.37	63.21	61.15	63.92	63.86	60.89	61.28	65.17	67.47	52.87	51.79	<u>65.18</u>	<u>70.10</u>	67.91	71.28
	Std	2.14	1.87	1.85	3.72	1.90	1.69	1.44	1.13	2.75	3.23	<u>1.02</u>	3.91	0.80	<u>1.28</u>	2.36	2.11
20	1	62.66	61.42	62.55	64.38	60.10	60.91	58.75	59.22	62.40	62.74	51.37	49.81	64.89	<u>68.32</u>	<u>64.02</u>	69.48
	2	61.58	63.38	61.22	63.30	61.96	64.00	58.25	60.10	<u>64.50</u>	67.31	51.44	58.02	64.08	<u>70.40</u>	65.19	72.94
	3	<u>63.82</u>	<u>64.27</u>	60.13	51.67	57.58	57.25	60.26	59.52	59.67	60.19	51.79	54.99	61.94	<u>70.16</u>	67.50	72.06
	4	<u>65.25</u>	67.91	61.12	62.23	60.86	61.45	61.83	63.15	62.31	66.04	50.78	62.28	63.91	<u>69.78</u>	68.29	73.56
	5	62.54	66.07	63.31	56.26	61.93	64.22	60.97	62.04	<u>64.07</u>	67.66	55.10	55.29	61.60	<u>68.45</u>	67.82	71.41
	6	63.99	67.09	<u>66.21</u>	62.59	63.54	65.45	60.00	61.58	65.25	68.32	52.86	57.22	65.84	<u>70.34</u>	67.57	72.08
	7	60.54	60.97	59.68	49.65	<u>60.96</u>	59.18	59.15	59.15	59.49	59.02	52.09	53.13	63.57	67.69	59.02	<u>62.68</u>
	8	61.73	63.72	63.79	66.42	63.21	63.28	60.85	60.33	67.15	66.27	53.79	56.57	<u>67.53</u>	<u>70.85</u>	71.01	71.59
	9	<u>64.50</u>	67.11	58.42	62.28	61.56	61.54	58.58	59.41	61.38	66.73	52.95	52.13	62.00	<u>68.15</u>	68.97	71.84
	10	67.65	68.27	65.62	59.87	66.24	65.18	60.37	60.77	<u>67.05</u>	69.07	52.91	56.96	64.35	<u>71.29</u>	64.41	71.72
	11	62.84	64.83	<u>63.91</u>	63.95	61.18	63.32	57.25	58.87	63.64	65.17	50.23	49.62	63.90	<u>69.48</u>	67.61	72.08
	12	64.32	67.87	61.96	62.81	63.14	63.62	60.68	60.79	<u>64.80</u>	66.22	53.72	58.49	64.28	<u>70.93</u>	70.00	73.76
	Avg	63.45	65.24	62.33	60.45	61.85	62.45	59.74	60.41	63.48	65.39	52.42	55.38	<u>63.99</u>	<u>69.65</u>	66.78	71.27
	Std	1.90	2.53	2.36	5.22	2.12	2.46	1.34	<u>1.32</u>	2.52	3.15	<u>1.40</u>	3.71	1.68	1.22	3.22	2.92
40	1	58.40	60.63	59.09	52.63	<u>61.46</u>	61.98	57.46	57.21	59.92	62.93	52.20	51.39	63.60	69.37	60.14	<u>65.90</u>
	2	55.72	59.84	57.33	45.74	59.51	62.50	56.46	58.85	60.47	62.10	51.53	50.27	<u>63.04</u>	<u>69.43</u>	64.61	71.91
	3	61.09	65.00	58.31	54.70	60.57	59.58	57.75	58.30	60.77	60.06	51.29	44.09	<u>62.83</u>	<u>69.12</u>	66.02	71.05
	4	62.62	67.18	59.44	63.46	63.02	64.25	57.60	59.23	<u>64.50</u>	68.56	51.83	58.48	61.79	<u>68.84</u>	65.22	70.68
	5	61.58	64.60	<u>64.05</u>	63.00	57.99	58.31	56.78	57.05	60.37	59.96	51.18	54.44	62.43	<u>68.49</u>	64.34	71.55
	6	60.53	64.20	<u>62.97</u>	59.95	60.78	61.72	57.77	58.43	<u>63.43</u>	66.76	52.83	53.18	61.78	<u>69.85</u>	64.61	72.75
	7	59.17	61.58	59.64	52.33	<u>61.44</u>	60.33	56.30	56.25	59.08	60.06	51.11	49.93	63.75	69.19	60.12	<u>66.69</u>
	8	65.24	67.27	60.86	60.76	<u>66.15</u>	63.91	59.00	58.26	67.88	68.32	53.46	53.20	66.11	70.49	65.07	<u>68.50</u>
	9	57.08	61.80	<u>62.72</u>	65.20	61.13	62.82	56.73	56.95	61.19	63.86	49.79	61.90	61.59	<u>68.42</u>	65.23	69.88
	10	<u>63.49</u>	63.19	62.13	55.84	60.08	61.04	56.81	55.78	60.22	62.01	52.11	49.37	67.77	70.16	61.96	<u>69.00</u>
	11	59.14	62.12	<u>61.70</u>	57.92	59.85	61.68	57.92	58.34	60.20	62.78	52.94	57.77	62.81	<u>68.36</u>	59.84	70.59
	12	62.54	64.21	59.35	56.73	<u>62.57</u>	62.93	60.10	59.81	63.10	63.96	52.98	48.00	67.49	<u>71.05</u>	<u>64.87</u>	72.02
	Avg	60.55	63.47	60.63	57.35	61.21	61.75	57.56	57.87	61.76	63.45	51.94	52.67	63.75	<u>69.40</u>	<u>63.50</u>	70.04
	Std	2.76	2.38	2.08	5.57	2.06	1.73	<u>1.10</u>	<u>1.22</u>	2.52	3.03	1.03	4.95	2.18	0.85	2.30	2.14

Table 5: Full results of the label disturbance experiment on **Sleep** data. The **best results** are in bold and we underline the **second best results**.

		Noisy Label Learning						Time Series Classification						Both																			
		SIGUA		UNICON		Sel-CL		MiniRocket		TimesNet		PatchTST		SREA		Con4m																	
r%	Exp	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁																
0	1	54.74	54.79	63.40	62.41	<u>63.86</u>	<u>63.49</u>	62.80	62.16	59.92	58.73	58.95	58.42	49.76	48.95	69.31	68.80																
	2	52.76	52.69	<u>63.15</u>	62.49	62.71	<u>62.87</u>	61.32	61.14	59.48	59.72	59.04	58.60	48.13	46.93	67.54	67.63																
	3	56.24	56.19	63.40	62.62	<u>65.73</u>	<u>65.47</u>	63.49	62.74	61.47	60.76	60.21	59.44	50.32	49.38	69.14	69.29																
	4	53.83	53.51	61.21	61.01	<u>62.72</u>	<u>62.88</u>	62.16	61.64	58.17	58.23	59.13	59.22	49.59	48.55	66.61	66.66																
	5	54.82	54.36	63.19	62.89	<u>63.62</u>	<u>63.82</u>	62.30	62.20	61.14	60.80	57.83	57.45	49.79	48.82	66.55	66.61																
	6	54.43	54.12	61.89	62.11	61.95	<u>62.37</u>	<u>62.07</u>	62.10	59.02	58.76	57.17	57.25	50.77	50.24	68.43	69.11																
	Avg	Std	54.47	1.15	54.28	1.19	62.71	0.93	62.26	0.66	<u>63.43</u>	1.32	<u>63.48</u>	1.10	62.36	0.73	62.00	0.55	59.87	1.26	59.50	1.10	58.72	1.07	58.40	0.90	49.73	0.90	48.81	1.09	67.93	1.22	68.02
20	1	54.24	53.73	63.38	62.75	<u>64.13</u>	<u>64.41</u>	62.30	61.86	59.58	58.07	57.14	56.82	50.00	49.80	67.57	67.07																
	2	51.72	51.04	63.01	62.68	<u>63.29</u>	<u>63.58</u>	61.91	61.51	59.84	57.44	56.12	55.50	48.17	47.56	64.01	64.25																
	3	54.68	54.51	62.44	61.44	<u>64.29</u>	<u>64.58</u>	62.95	62.35	59.51	56.10	57.53	57.03	50.63	49.30	68.76	68.50																
	4	53.33	53.12	61.25	59.39	<u>62.34</u>	<u>62.33</u>	62.28	61.61	57.14	57.23	57.32	56.77	48.82	47.65	65.57	65.25																
	5	53.20	52.83	62.72	61.92	<u>63.15</u>	<u>63.28</u>	61.90	61.75	60.00	58.99	55.18	54.78	48.48	48.18	66.26	65.90																
	6	53.82	53.18	<u>62.73</u>	61.60	61.97	<u>62.51</u>	61.69	61.43	58.97	58.47	56.85	56.05	50.45	50.28	67.49	66.86																
	Avg	Std	53.50	1.03	53.07	1.16	62.59	0.73	61.63	1.22	<u>63.19</u>	0.93	<u>63.45</u>	0.94	62.17	0.45	61.75	0.33	59.17	1.06	57.72	1.02	56.69	0.89	56.16	0.89	49.43	1.06	48.80	1.16	66.61	1.69	66.31
40	1	53.08	52.10	60.95	58.17	<u>61.83</u>	<u>61.54</u>	59.57	58.62	57.61	57.20	56.78	55.98	48.99	47.23	66.79	65.38																
	2	51.21	50.08	60.47	58.12	<u>61.58</u>	<u>61.64</u>	58.62	57.96	56.62	55.26	53.94	52.60	46.15	44.56	65.60	64.27																
	3	54.12	53.85	61.02	59.63	<u>63.70</u>	<u>63.27</u>	60.03	59.18	55.81	55.30	53.72	52.12	48.97	47.98	66.31	65.36																
	4	52.38	52.21	60.43	57.58	<u>61.80</u>	<u>61.59</u>	59.41	58.68	55.38	54.36	55.06	54.31	48.10	45.53	66.02	65.69																
	5	50.99	49.48	59.88	57.16	<u>61.44</u>	<u>61.28</u>	58.45	57.47	57.43	56.80	52.47	50.80	48.39	44.03	63.54	61.82																
	6	51.19	50.18	<u>61.13</u>	59.40	60.76	<u>61.01</u>	59.06	58.36	57.24	55.44	53.32	52.50	48.71	45.00	63.76	63.33																
	Avg	Std	52.16	1.26	51.32	1.68	60.65	0.48	58.34	0.99	<u>61.85</u>	0.98	<u>61.72</u>	0.80	59.19	0.60	58.38	0.60	56.68	0.92	55.73	1.06	54.21	1.51	53.05	1.82	48.22	1.07	45.72	1.56	65.34	1.36	64.31

Table 6: Full results of the label disturbance experiment on **SEEG** data. The **best results** are in bold and we underline the **second best results**.

		Noisy Label Learning						Time Series Classification						Both			
		SIGUA		UNICON		Sel-CL		MiniRocket		TimesNet		PatchTST		SREA		Con4m	
r%	Exp	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁	Acc.	F ₁
SEEG	1	66.26	52.27	<u>69.46</u>	60.64	68.62	<u>62.57</u>	68.96	62.11	66.13	50.25	65.80	58.51	65.90	53.09	74.70	72.26
	2	67.93	55.21	<u>71.09</u>	<u>64.43</u>	64.73	51.65	69.02	61.12	65.70	50.55	68.12	60.60	64.29	57.57	75.23	73.21
	3	66.40	52.09	67.11	56.51	<u>72.02</u>	<u>67.27</u>	68.39	63.92	66.24	52.17	65.86	56.22	65.15	54.99	73.87	70.52
	Avg	66.87	53.19	<u>69.22</u>	60.53	68.46	60.50	68.79	<u>62.39</u>	66.02	50.99	66.59	58.45	65.11	55.21	74.60	72.00
	Std	0.93	1.75	2.00	3.96	3.65	8.01	0.35	1.42	0.28	1.03	1.32	2.19	0.81	2.25	0.69	<u>1.36</u>

689 H CASE STUDY

690 As shown in Figure 7, we present four cases to compare and demonstrate the differences between
691 our proposed *Con4m* and other baselines. The first two cases involve transitions from a seizure state
692 of label 1 to a normal state of label 0. The third case consists of entirely normal segments, while
693 the fourth case comprises entirely seizure segments. As illustrated in the figure, *Con4m* exhibits
694 more coherent narratives by constraining the predictions to align with the contextual information of
695 the data. Moreover, it demonstrates improved accuracy in identifying the boundaries of transition
696 states. In contrast, MiniRocket and SREA exhibit fragmented and erroneous predictions along the
697 time segments. This verifies that *Con4m* can achieve clearer recognition of boundaries, and it can
698 also make better predictions on the continuous time segments belonging to the same class.

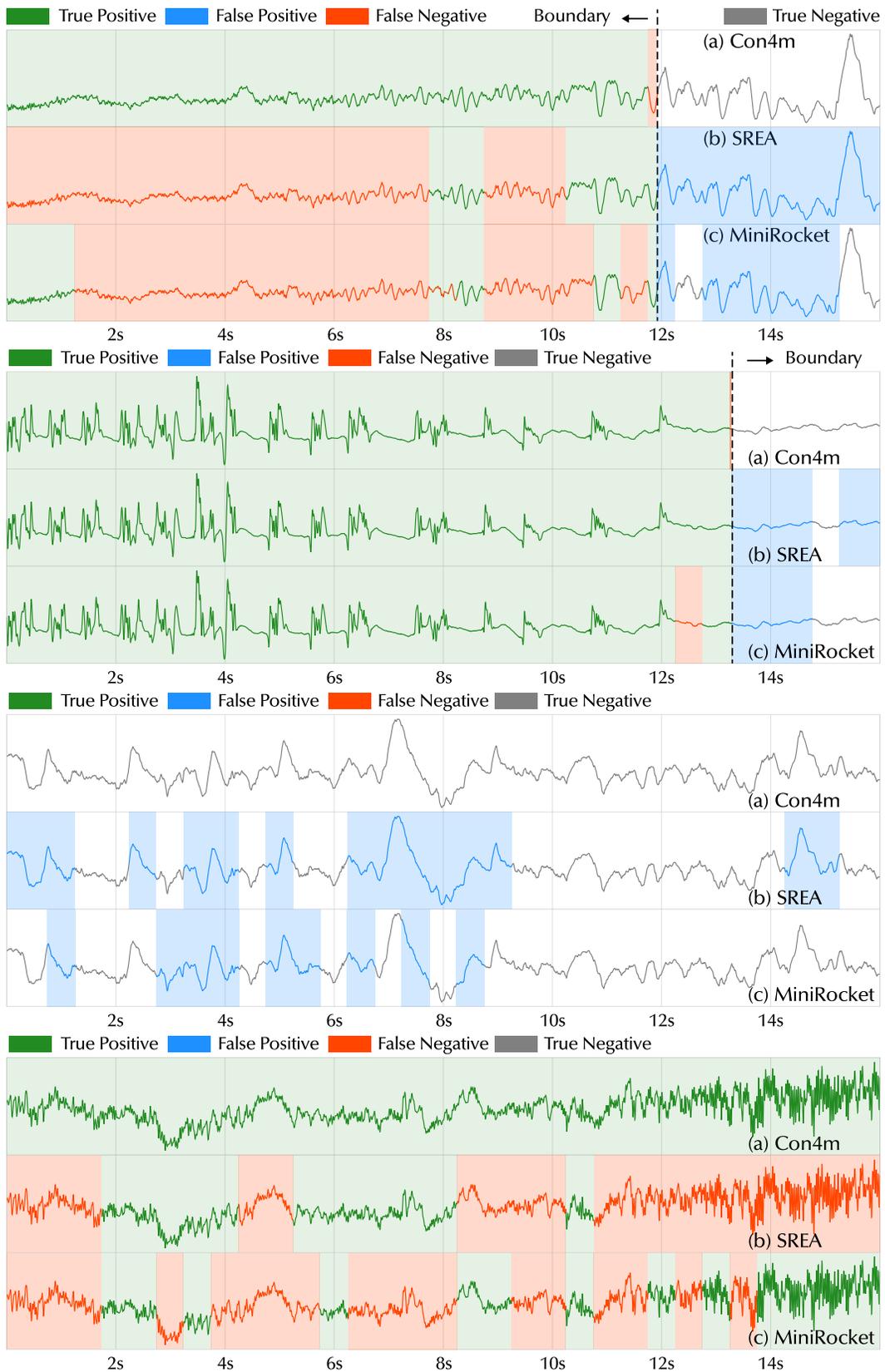


Figure 7: More cases for continuous time intervals in SEEG testing set.

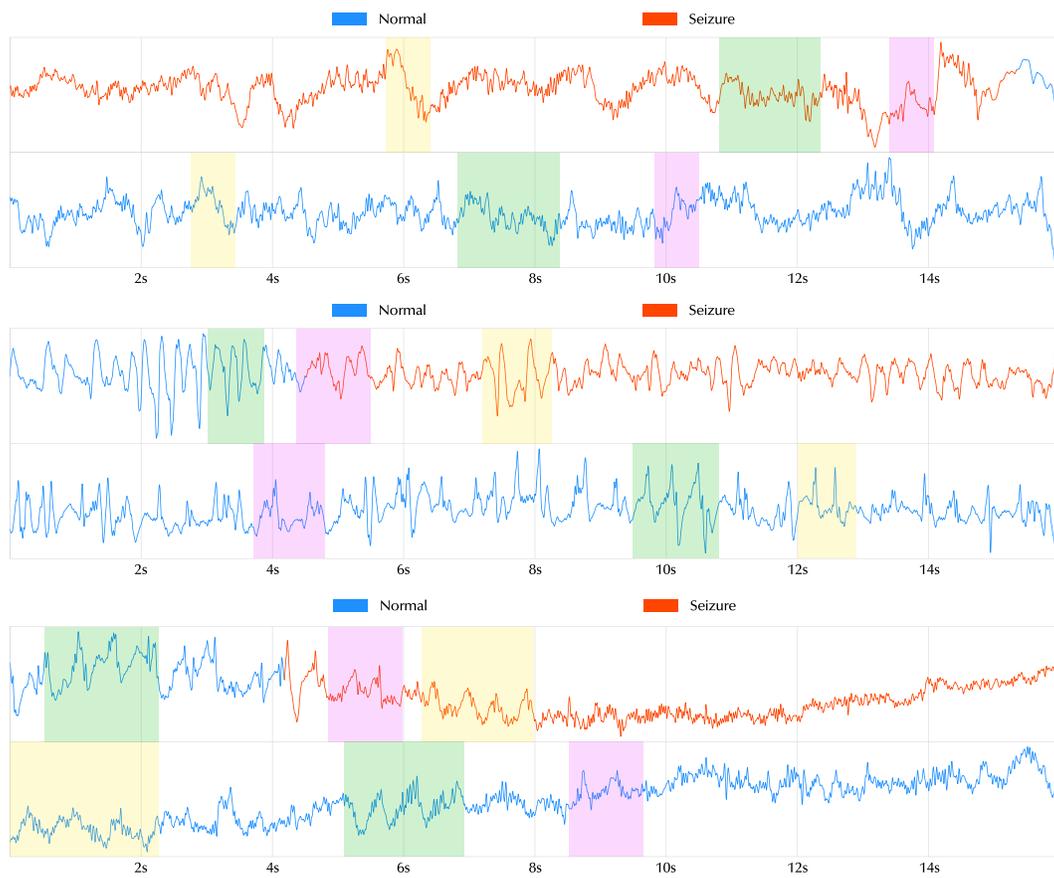


Figure 8: Comparison between boundary and central time intervals in SEEG data.