
RCT-MNM: Region-Channel-Temporal Masked Neural Modeling for Emotion and Psychiatric Disorder Analysis

Anonymous Authors¹

Abstract

We propose Region-Channel-Temporal Masked Neural Modeling (RCT-MNM), a self-supervised learning framework for EEG and sEEG-based emotion recognition and psychiatric disorder analysis. RCT-MNM integrates region-channel-temporal self-attention to model spatial-temporal dependencies and conditioned self-supervised learning for robust representation alignment. A post-pretraining refinement step further enhances task performance. Experimental results on multiple EEG and sEEG datasets demonstrate state-of-the-art performance in emotion recognition and task engagement classification. Additionally, RCT-MNM provides region- and channel-level interpretability, revealing meaningful neural dynamics for cognitive health applications. Our framework advances neural representation learning, enhancing the generalizability and interpretability of EEG/sEEG-based models for brain signal decoding and therapeutic insights.

1. Introduction

Psychiatric and mood disorders such as depression and anxiety, along with cognitive impairments, impact millions of individuals worldwide, significantly affecting mental well-being and cognitive function (Organization, 2001). These conditions are often difficult to diagnose and treat due to their complex neurophysiological underpinnings and reliance on subjective assessments (Kret & Ploeger, 2015; Goschke, 2014). As a result, there is a growing interest in leveraging neurophysiological signals, to develop objective and data-driven diagnostic tools for emotion recognition and cognitive state analysis (Houssein et al., 2022). Understanding the neural basis of these disorders can lead to improved treatment strategies, brain-computer interfaces

(BCIs), and personalized mental health interventions (Ng & Weisz, 2016).

Emotion recognition and cognitive analysis from electroencephalography (EEG) and stereo-EEG (sEEG) are essential for understanding human affective states and psychiatric disorders (Drane et al., 2021). However, these neural signals exhibit complex spatio-temporal dependencies, posing significant challenges for machine learning models in achieving robust and accurate classification (Gao et al., 2019). Traditional approaches rely on handcrafted features or shallow models, which fail to capture the intricate interactions across spatial and temporal scales (Rouast et al., 2019; Liu et al., 2024). Recent advancements in deep learning, particularly transformer-based models, have shown promise in neural representation learning in global level (Vaswani et al., 2017; Dosovitskiy et al., 2021; He et al., 2022; Assran et al., 2023).

Despite advancements in EEG/sEEG neural modeling, several challenges remain. (1) Hierarchical Spatial Representation: Existing models overlook the structured organization from electrodes to brain regions, limiting spatial feature extraction (Lachaux et al., 2003; Gevins et al., 1995). (2) Long-Term Temporal Dependencies: Conventional models struggle to capture sustained neural patterns essential for emotion and cognitive analysis. (3) Structured Attention Mechanisms: Global self-attention for spatial representation learning is computationally expensive, requiring more efficient structured attention. (4) Generalization Across Datasets: High inter-subject variability hinders model robustness (Saha & Baumert, 2020). Self-supervised learning (SSL) can improve feature extraction without labeled data (Rafiei et al., 2024; He et al., 2022; Selvaraju et al., 2017; Jiang et al., 2024; Wang et al., 2024). Addressing these challenges requires a model that integrates spatio-temporal dependencies, regional interactions, and efficient SSL strategies for improved EEG-based emotion and cognitive analysis.

We introduce **Region-Channel-Temporal Masked Neural Modeling (RCT-MNM)**, a self-supervised learning framework that leverages hierarchical attention mechanisms to capture multi-scale dependencies in EEG and sEEG signals. Our contributions are as follows:

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

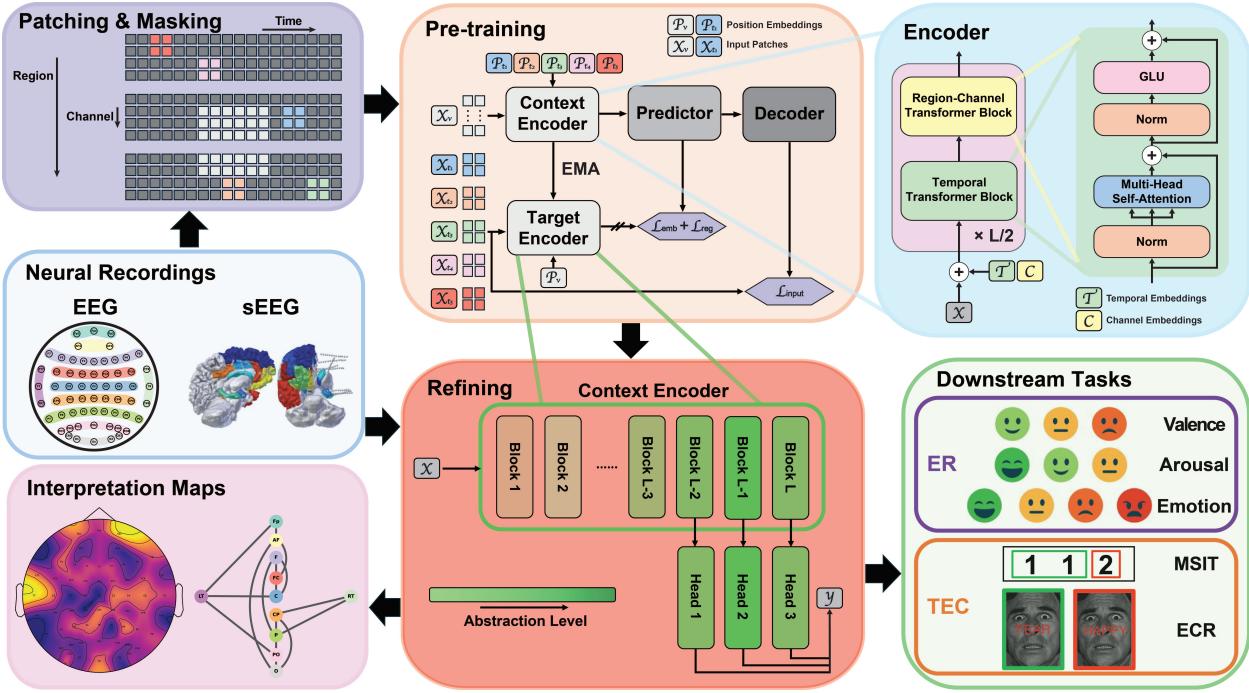


Figure 1. **RCT-MNM.** The model consists of three stages: Patching & Masking, segmenting and masking EEG/sEEG signals; Pre-training, learning region-channel-temporal representations; and Refining, fine-tuning the Encoder for improved downstream performance. It supports emotion recognition (ER) and task engagement classification (TEC) with interpretable neural representations.

1. We propose a novel Temporal-Region-Channel encoder that captures hierarchical neural dependencies in EEG and sEEG data through a structured transformer-based approach.
2. We introduce a Region-Channel Self-Attention Mechanism (RCSA) mechanism that models functional interactions between brain regions and electrode channels, enabling structured spatial dependency modeling.
3. We develop a region-conditioned self-supervised learning paradigm that employs Region-Channel-Temporal masking for representation alignment and input reconstruction to enhance model generalizability across diverse EEG/sEEG tasks.
4. We introduce a post-pre-training Adaptive Attention-Based Layer Refinement step that fine-tunes intermediate representations using to improve task-specific classification performance.
5. Our model achieves state-of-the-art classification performance for emotion recognition and task engagement classification in subject-dependent and subject-independent settings.
6. We enhance interpretability by generating region- and channel-level activation maps and self-attention graphs,

providing insights into functional brain dynamics for affective and psychiatric disorder analysis.

2. Related Work

Graph-Based Models for EEG Analysis. Graph-based models have been widely explored for EEG-based emotion recognition and cognitive state classification. Dynamic Graph Convolutional Neural Networks (DGCNN) dynamically construct EEG graph representations and apply graph convolution to model spatial dependencies among electrode channels (Song et al., 2018). Regularized Graph Neural Networks (RGNN) integrate a neuroscience-inspired adjacency matrix and employ node-wise domainadversarial training (NodeDAT) for cross-subject adaptation, and emotion-aware distribution learning (EmotionDL) for noise robustness in EEG-based emotion recognition. (Zhong et al., 2020). Local-Global Graph Networks (LGG-Net) extend this approach by leveraging local and global graph representations, enhancing the model’s ability to capture fine-grained spatial structures (Ding et al., 2023). HDGCN employs dual-branch hierarchical dynamic graph convolution to capture global and local EEG spatial patterns, integrating a layerwise adjacency matrix and an auxiliary information module for enhanced feature fusion and interpretability (Ye et al., 2022). However, these models exhibit limited

110 spatio-temporal integration and generalizability in spatial
 111 representations.

112 **Temporal Modeling for Neural Signals.** Temporal de-
 113 pendencies in neural signals have been modeled using con-
 114 volutional and transformer-based architectures. TSception
 115 employs multi-scale 1D convolutional kernels to capture
 116 both short- and long-term temporal dependencies while
 117 leveraging spatial asymmetry in EEG signals (Ding et al.,
 118 2022). EEG Conformer introduces a hybrid architecture
 119 combining convolution and self-attention to capture local
 120 and global temporal dependencies in EEG data (Song et al.,
 121 2022). Despite their effectiveness, convolution-based meth-
 122 ods struggle with long-range dependencies, and temporal
 123 self-attention itself does not explicitly encode inductive bi-
 124 ases about EEG/sEEG signal structures, which can make
 125 training inefficient and lead to overfitting on small datasets.
 126

127 **Transformer-Based Self-Supervised Learning for EEG.**
 128 As self-supervised learning has gained traction in computer
 129 vision, there is increasing interest in adapting these tech-
 130 niques to the neural domain, where the inherent spatio-
 131 temporal structure of EEG signals presents unique chal-
 132 lenges for representation learning. Masked Autoencoders
 133 (MAEEG) perform masked input reconstruction to learn ro-
 134 bust EEG representations (Chien et al., 2022). MMM unifies
 135 EEG channel topologies and employs multi-dimensional po-
 136 sition encoding, hierarchical channel representations, and
 137 multi-stage pre-training to enhance cross-dataset generaliza-
 138 tion for EEG-based emotion recognition (Yi et al., 2024).
 139 LaBraM introduces a biologically inspired approach to self-
 140 supervised EEG representation learning, leveraging latent
 141 brain dynamics modeling to capture underlying neural pro-
 142 cesses across diverse cognitive states (Jiang et al., 2024).
 143 EEGPT extends large-scale pre-training techniques to EEG
 144 analysis, leveraging a self-supervised dual pretraining strat-
 145 egy that combines spatio-temporal representation alignment
 146 and masked reconstruction (Wang et al., 2024). SEEG-
 147 nificant framework tokenizes sEEG signals with convolu-
 148 tions, applies temporal and spatial self-attention, integrates
 149 electrode locations, and employs subject-specific heads for
 150 multi-subject behavioral decoding and cross-subject transfer
 151 learning (Mentzelopoulos et al., 2024). However, they either
 152 lack hierarchical spatial modeling or fail to effectively inte-
 153 grate long-range temporal dependencies in EEG and sEEG
 154 signals.

155 3. Methodology

156 3.1. Problem Setting and Formulation

157 We represent the neural data (EEG/sEEG) as $\mathbf{X}_{\text{input}} \in$
 158 $\mathbb{R}^{C \times T}$, where C is the number of channels, and T is the
 159 number of time points. The problem of task engagement
 160 classification (TEC) and low-high valence/arousal classi-
 161

162 fication is formulated as binary classification tasks pre-
 163 dicting rest-stage versus task-stage and low versus high
 164 valence/arousal with labels $y \in \{0, 1\}$. The problem of
 165 valence recognition is formulated as a multi-class classifi-
 166 cation task predicting different types of emotions with labels
 167 $y \in \{0, \dots, E - 1\}$, where E is the number of types of emo-
 168 tions. The goal for these tasks is to predict the label y based
 169 on a given neural input $\mathbf{X}_{\text{input}}$.

170 We first segment the data into non-overlapping patches.
 171 The input patch representation is defined as: $\mathbf{X}_{\text{input}} =$
 172 $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^{1 \times P}$, where P is the tempo-
 173 ral patch size and $N = \frac{CT}{P}$ is the number of patches.
 174 Each patch undergoes a linear projection, such that $\mathbf{X} =$
 175 $\mathbf{X}_{\text{input}} W_P$, where $\mathbf{X} \in \mathbb{R}^{N \times d}$, and d is the embedding size.
 176

177 3.2. Temporal-Region-Channel Encoder

178 The Temporal-Region-Channel Encoder is designed to cap-
 179 ture the intricate spatio-temporal dependencies in EEG and
 180 sEEG data through a hierarchical transformer-based archi-
 181 tecture. It consists of two primary components: the Region-
 182 Channel Transformer Block and the Temporal Transformer
 183 Block, which work together to extract generalized neural
 184 representations.

185 The Region-Channel Transformer Block models spatial de-
 186 pendencies across EEG channels by applying a multi-head
 187 self-attention mechanism across the region-channel space.
 188 To construct region-level representations, we add extra re-
 189 gion tokens to the channels tokens and use Region-Channel
 190 Self-Attention mechanism (**RCSA**) to learn region represen-
 191 tations from channel-level information based on pre-defined
 192 functional region-channel map and region-channel attention
 193 mask. The region-channel maps and corresponding atten-
 194 tion masks are defined based on the location of electrodes
 195 on function areas for EEG (Alarcao & Fonseca, 2017), and
 196 using electrode labeling algorithm (ELA) for sEEG (Peled
 197 et al., 2017). RCSA is a pivotal element in capturing both
 198 localized and global interactions within neural data, which
 199 especially captures the functional connectivity and coordi-
 200 nation among neural patches from distinct brain regions.
 201 These regional-channel dynamics are essential for under-
 202 standing higher-order brain functions, such as task execu-
 203 tion or cognitive processing, which involve communication
 204 across spatially distributed networks.

205 Given an input representation \mathbf{X} consisting of channel to-
 206 kens $\mathbf{X}_c \in \mathbb{R}^{C \times d}$, we append randomly initialized re-
 207 gion tokens $\mathbf{X}_r \in \mathbb{R}^{R \times d}$, forming the combined input:
 208 $\mathbf{X}_{cr} = [\mathbf{X}_c; \mathbf{X}_r] \in \mathbb{R}^{(C+R) \times d}$, where R is the number of
 209 regions. We define a binary region-channel attention mask
 210 $M_{rc} \in \{0, 1\}^{(C+R) \times (C+R)}$ based on the functional region-
 211 channel map, ensuring that each region token only attends to
 212 its corresponding channels and other regions, and each chan-
 213 nel only attends its corresponding region and other channels

within the same region. The region-channel self-attention mechanism (**RCSA**) operates as follows:

$$\begin{aligned} \mathbf{Q} &= \mathbf{X}_{cr} W_Q, \quad \mathbf{K} = \mathbf{X}_{cr} W_K, \quad \mathbf{V} = \mathbf{X}_{cr} W_V \\ A_{rc} &= \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \odot M_{rc} \right) \quad (1) \\ \mathbf{X}'_{cr} &= A_{rc} \mathbf{V} \end{aligned}$$

where W_Q, W_K, W_V are linear projection matrices for Query, Key, and Value, and A_{rc} represents the region-channel attention matrix, capturing dependencies across neural channels and regions, and \odot denotes element-wise multiplication.

These region-channel tokens further undergo an additional feedforward network to enhance spatial modeling (Shazeer, 2020):

$$\begin{aligned} \mathbf{X}''_{cr} &= \text{FFN}_{\text{SwiGLU}}(\mathbf{X}'_{cr}) \\ &= (\text{Swish}_1(\mathbf{X}'_{cr} W) \odot \mathbf{X}'_{cr} V) W_2 \quad (2) \end{aligned}$$

where $\text{Swish}_\beta(x) = x\sigma(\beta x)$. By the adding residual connections, the Region-Channel Transformer Block can be represented as:

$$\begin{aligned} \mathbf{X}'_{cr} &= \text{RCSA}(\text{Norm}(\mathbf{X}_{cr})) + \mathbf{X}_{cr} \\ \mathbf{X}''_{cr} &= \text{FFN}_{\text{SwiGLU}}(\text{Norm}(\mathbf{X}'_{cr})) + \mathbf{X}'_{cr} \quad (3) \end{aligned}$$

To model long-range dependencies in neural activity, we apply a Temporal Transformer Block, which operates across the time dimension while preserving spatial information. The attention mechanism is a standard temporal self-attention (TSA), and the same feedforward network is used for feature mixing. Therefore, the Temporal Transformer Block can be formulated as:

$$\begin{aligned} \mathbf{X}'_t &= \text{TSA}(\text{Norm}(\mathbf{X}_t)) + \mathbf{X}_t \\ \mathbf{X}''_t &= \text{FFN}_{\text{SwiGLU}}(\text{Norm}(\mathbf{X}'_t)) + \mathbf{X}'_t \quad (4) \end{aligned}$$

where $\mathbf{X}_t \in \mathbb{R}^{S \times d}$, and $S = \frac{T}{P}$ is the number of time segments.

We construct the Encoder in a two-stage hierarchical manner. First, we apply Temporal Transformer Block to model long-range dependencies across time, ensuring that dynamic temporal patterns are well-represented. Then, we employ the Region-Channel Transformer Block to integrate spatial dependencies within and across functional brain regions. This design choice is motivated by the fundamental structure of neural signals, where temporal dynamics exhibit continuity and correlation over time, while spatial interactions are constrained by anatomical and functional organization. By first encoding temporal dependencies, we allow the model to learn stable temporal representations before aggregating region-level spatial information. This sequential processing also improves computational efficiency by reducing spatial

redundancy in raw channel-wise inputs. The Predictor is also in the same structure as the Encoder, and the Decoder is constructed in a Region-Channel-then-Temporal Transformer Block manner for reversed reconstruction.

To enhance the representation of temporal and spatial dependencies, we incorporate Rotary Position Embeddings (RoPE) (Su et al., 2024) for the temporal dimension and learnable spatial embeddings (Dosovitskiy et al., 2021) for both channel and region representations. This design choice is driven by the strong inductive bias introduced by RCSA, which enforces structured connectivity constraints between channels and brain regions.

3.3. Region-Conditioned Masked Neural Modeling

Our proposed framework utilizes Masked Neural Modeling (MNM) as an input-representation predictive architecture for self-supervised learning of EEG and sEEG patches. The model consists of three core components: a context encoder $f(\cdot; \theta)$, a target encoder $f(\cdot; \tilde{\theta})$, and a predictor $g(\cdot; \psi)$, which collectively learn to align target representations. Additionally, a decoder $h(\cdot; \phi)$ is applied for input reconstruction, ensuring robust feature generalizations across spatio-temporal dimensions. We further introduce details about our proposed masking strategy **RCT-Mask**, **Regional Conditioning** for guiding the learning process, and **MNM-Loss** function to improve generalization of neural representations.

3.3.1. REGION-CHANNEL-TEMPORAL MASKING

Our pre-training framework is designed to predict neural inputs \mathbf{X} and representations \mathbf{Z} across multiple target blocks based on the representation of a context block. In training, the model learns to predict representations of other sections of the neural data within the latent space based on a single observation. However, random sampling of targets, as seen in methods like MAE, may lead the model to rely on simplistic shortcuts, such as interpolating time series or defaulting to more prevalent patterns in the data, potentially reducing its generalizability (Assran et al., 2023). It is important to consider that neural patches exhibit spatial variability based on functional brain organization, as well as temporal variability related to brain states and task conditions. Additionally, complex nonlinear relationships among brain networks add layers of interaction across different brain patches, further challenging the model's capacity to learn robust, generalized representations.

The Region-Channel-Temporal Masking (**RCT-Mask**) approach is a foundational component of our pre-training framework. It is designed to encourage the model to learn robust and generalized representations from neural data by masking signals selectively across temporal, and localized/distributed spatial domains, either independently or in

combination. This mechanism prevents the model from relying on simplistic patterns or interpolation, driving it to uncover deeper, nonlinear relationships.

For a given EEG/sEEG input $\mathbf{X} \in \mathbb{R}^{N \times d}$, where $N = C \times S$, and C is the number of channels, and S is the number of time sequences, the framework identifies a context block \mathbf{x}_c . This context block is randomly selected from the input neural segments within the range $\{\rho_c, \eta_c\}$, where ρ_c and η_c are the range ratio for the context block along the spatial and temporal dimensions respectively. The corresponding binary mask for the context block can be defined as $\mathbf{M}_c \in \{0, 1\}^{C \times S}$. The context block serves as the observation basis for predicting masked neural representations of other segments, termed target blocks \mathbf{x}_t , with $\mathbf{x}_c, \mathbf{x}_t \in \mathbf{X}$ and $\mathbf{x}_c \cap \mathbf{x}_t = \emptyset$. The objective of the model is to predict \mathbf{x}_t and \mathbf{z}_t from \mathbf{x}_c , where \mathbf{z} represents the latent representations derived from the encoders, and \mathbf{z}_t is the latent representations of the target blocks.

We categorize the remaining patches excluding the context block into five distinct and non-overlapping domains: cross-time, cross-channel, cross-time-channel, cross-region, cross-region-channel. We randomly sample multiple blocks from each of the five domains as targets blocks, with binary masks $\mathbf{M} \in \{0, 1\}^{C \times S}$ corresponding to each domain as $\mathbf{M}^t, \mathbf{M}^c, \mathbf{M}^{ct}, \mathbf{M}^r$, and \mathbf{M}^{rt} , within the range $\{\rho_t, \eta_t\}$, where ρ_t and η_t are the range ratio for the target blocks along the spatial and temporal dimensions respectively. Therefore, the RCT-Mask \mathbf{M}_{RCT} , context block, and target blocks can be written as:

$$\begin{aligned} \mathbf{M}_{RCT} &= \mathbf{M}^t \cup \mathbf{M}^c \cup \mathbf{M}^{ct} \cup \mathbf{M}^r \cup \mathbf{M}^{rt} \\ \mathbf{x}_c &= \mathbf{X} \odot \mathbf{M}_c \\ \mathbf{x}_t &= \mathbf{X} \odot \mathbf{M}_{RCT} \end{aligned} \quad (5)$$

3.3.2. REGIONAL CONDITIONING

Given a tokenized input representation $\mathbf{x} \in \mathbb{R}^{N \times d}$ and corresponding positional embeddings $\mathbf{p} \in \mathbb{R}^{N \times d}$, where N is the number of tokens, we define a self-supervised learning framework based on context-target prediction and reconstruction. Let c denote the set of indices corresponding to the context block, such that the context tokens and positional embeddings are represented as \mathbf{x}_c and \mathbf{p}_c . Similarly, let t represents target blocks with target tokens and positional embeddings \mathbf{x}_t and \mathbf{p}_t . The encoder function $f(\cdot; \theta)$ maps the context input to latent representations: $\mathbf{z}_c = f(\mathbf{x}_c; \theta)$, where θ denotes the encoder parameters. The learned context representations \mathbf{z}_c are then used to predict the representations of the masked target tokens $\mathbf{z}_t = f(\mathbf{x}; \tilde{\theta})$, where $\tilde{\theta}$ represents an exponential moving average (EMA) of the encoder weights θ to stabilize training. A predictor function $g(\cdot; \psi)$ further refines the predictions by incorporating both the context representations \mathbf{z}_c and the positional embeddings of the target tokens \mathbf{p}_t , where ψ are the learnable parameters of the predictor function. The decoder function $h(\cdot; \phi)$ finally reconstruct the input target tokens using predicted target representations $\hat{\mathbf{z}}_t$ and target positional embeddings, such that:

$$\hat{\mathbf{z}}_t = g(\mathbf{z}_c, \mathbf{p}_t; \psi), \quad \hat{\mathbf{x}}_t = h(\hat{\mathbf{z}}_t, \mathbf{p}_t; \phi) \quad (6)$$

where $\hat{\mathbf{z}}_t$ is the predicted target representation, $\hat{\mathbf{x}}_t$ is the reconstructed target input, and ϕ are the learnable parameters of the decoder function. The encoder function f , predictor function g , and decoder function h are implemented as transformer blocks.

Our proposed methodology also incorporates explicit regional conditioning into the encoder during pre-training and inference to enhance the self-supervised learning process. Specifically, the regional conditioning is achieved by appending region-specific positional embeddings as additional tokens in the input sequence processed by the transformer blocks. This formulation can be expressed as follows:

$$\mathbf{z}_c^t = f(\mathbf{x}_c, \mathbf{p}_t^{r_t}; \theta), \quad \mathbf{z}_t^c = f(\mathbf{x}, \mathbf{p}_c^{r_c}; \tilde{\theta}) \quad (7)$$

where \mathbf{z}_c^t represents the context embedding conditioned on other regional targets, while \mathbf{z}_t^c denotes the target representation conditioned on the regional contexts, and $\mathbf{p}_t^{r_t}$ and $\mathbf{p}_c^{r_c}$ are the region positional embedding for regions r_t and r_c , which are the unique regions existed in target and context blocks. Therefore, the each region positional embedding is only used once as an additional token, and all unique tokens are appended to the input tokens.

EEG and sEEG datasets often exhibit variations in the number of channels, sampling rates, and recording durations. Explicit temporal and channel-wise conditioning may require recalibration for different settings and reduce generalizability. By conditioning only on region-level embeddings, our model maintains flexibility across diverse datasets without requiring channel- and temporal-level adjustments. Furthermore, since each region embedding is only used once as an additional token in the input sequence, rather than being repeatedly appended per time segment or channel, the model avoids excessive computational overhead. This allows for faster training and inference while retaining critical distributed spatial dependencies.

3.3.3. MASKED NEURAL MODELING LOSS

The loss function for our framework (**MNM-Loss**) is a combination of the input reconstruction loss, the representation alignment loss, and the VICReg-inspired regularization terms, structured as follows:

$$\mathcal{L}_{MNM} = \lambda_1 \mathcal{L}_{input} + \lambda_2 \mathcal{L}_{rep} + \lambda_3 \mathcal{L}_{reg} \quad (8)$$

Input Reconstruction & Representation Alignment Loss. The MNM objective is optimized by the average l_2 distance

275 between the predicted patch-level representations/inputs of
 276 the target blocks $\{\hat{\mathbf{z}}_t, \hat{\mathbf{x}}_t\}$ and the ground-truth target representations/inputs $\{\mathbf{z}_t, \mathbf{x}_t\}$, which is formulated as the Input
 277 Reconstruction and Representation Alignment Loss:
 278

$$\mathcal{L}_{\text{input}} = \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 \quad \mathcal{L}_{\text{rep}} = \|\hat{\mathbf{z}}_t - \mathbf{z}_t\|_2^2 \quad (9)$$

281 **VIC Regularization.** To further improve representation
 282 learning, we incorporate Variance-Invariance-Covariance
 283 Regularization (VICReg) (Bardes et al., 2022), which
 284 encourages decorrelated feature representations, preventing
 285 collapse into trivial solutions. VICReg introduces two key
 286 constraints: (1) Variance Regularization is a hinge function
 287 on the standard deviation of the representation \mathbf{z} along the
 288 batch dimension, ensuring that all dimensions of the learned
 289 embedding space exhibit meaningful variation, enhancing
 290 feature diversity. (2) Covariance Regularization minimizes
 291 off-diagonal elements of the covariance matrix along the
 292 batch dimension, promoting statistical independence
 293 between learned features. It can be represented in a batch-wise
 294 format as:

$$\begin{aligned} \mathcal{L}_{\text{reg}} &= \lambda_{\text{var}} * \mathcal{L}_{\text{var}} + \lambda_{\text{cov}} * \mathcal{L}_{\text{cov}} \\ \mathcal{L}_{\text{var}} &= \frac{1}{d} \sum_{j=1}^d \max(0, \gamma - \sqrt{\text{Var}(\mathbf{z}^j) + \epsilon}) \\ \mathcal{L}_{\text{cov}} &= \frac{1}{d} \sum_{i \neq j} \left[\frac{1}{n-1} \sum_{i=1}^n (\mathbf{z}_i - \bar{\mathbf{z}})(\mathbf{z}_i - \bar{\mathbf{z}})^T \right]_{i,j}^2 \end{aligned} \quad (10)$$

304 where $\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$. $\mathbf{z}_i \in \mathbb{R}^d$ is the i^{th} vector in this
 305 batch, and \mathbf{z}^j is the vector composed of each value at
 306 dimension j in all vectors in the batch. By enforcing these
 307 constraints, our method ensures that learned features remain
 308 informative, diverse, and uncorrelated, ultimately improving
 309 the generalization of the learned neural representations.

3.4. Adaptive Attention-Based Layer Refinement

310 The last layer of the pre-trained encoder might achieve low
 311 representation quality because of splitting the encoding
 312 capability for predicting and decoding. To enhance the effec-
 313 tiveness of our pre-trained encoder for downstream tasks,
 314 we introduce a post-pre-training refinement process, which
 315 fine-tunes intermediate representations rather than relying
 316 solely on the final encoder layer. This approach leverages an
 317 ensemble of multi-layer perceptron (MLP) heads attached
 318 to multiple layers near the end of the encoder, allowing the
 319 model to extract more abstract feature representations. By
 320 distributing refinement across multiple layers, the model
 321 improves adaptability and generalization in downstream
 322 applications. Unlike conventional methods that depend ex-
 323 clusively on the last encoder layer for feature extraction, our
 324 refinement strategy aggregates embeddings from multiple
 325 late-stage layers, mitigating potential limitations of a single
 326 late layer representation.

Given an input \mathbf{x} , let $f_l(\cdot; \theta)$ be the function of l consecutive
 blocks in the encoder, where the layer-wise representations
 are computed as:

$$\mathbf{z}_l = f_l(\mathbf{x}; \theta), \quad l \in \{1, \dots, L\} \quad (11)$$

where θ represents the encoder weights, and L is the total
 number of encoder layers. To refine the feature representa-
 tions for downstream tasks, we aggregate representations
 from the final k layers of the encoder. We introduce **Adap-
 tive Attention-Based Layer Refinement** to dynamically
 select important layers. For each layer, we can compute its
 attention score a_l as:

$$a_l = \text{softmax} \left(\frac{\mathbf{q}^\top W_q \text{Norm}(\mathbf{z}_l)}{\sqrt{d}} \right) \quad (12)$$

where W_q is a learnable query projection matrix, and \mathbf{q} is a
 global learnable query vector. The final refined embedding
 \mathbf{z}_{ref} is then adaptively aggregated over k layers as:

$$\mathbf{z}_{\text{ref}} = \sum_{l=L-k+1}^L \alpha_l \text{Norm}(\mathbf{z}_l) \quad (13)$$

which can be further attached to a classification head for
 downstream tasks. This adaptive refinement mechanism en-
 sures that informative layers are prioritized while redundant
 layers are suppressed. The resulting refined embeddings
 provide a more stable and contextually meaningful repre-
 sentation for downstream classification, ensuring improved
 performance across neural cognitive and affective tasks.

4. Experiments

4.1. Experimental Setup

Datasets. We used three publicly available datasets for
 the EEG emotion recognition (ER) tasks: **SEED** (Zheng &
 Lu, 2015) for three-class valence classification; **SEED-IV**
 (Zheng et al., 2018) for four-class valence classification; and
DEAP (Koelstra et al., 2011) for low-high valence/arousal
 binary classification. We used two datasets MSIT and ECR
 (Provenza et al., 2019) for the sEEG task engagement classi-
 fication (TEC). We employed Subject-Dependent (SD) and
 Subject-Independent (SI) paradigms for SEED, SEED-IV,
 and DEAP, and we implemented SD paradigm for MSIT
 and ECR. A detailed description of datasets can be found in
 Appendix B.

Baselines. To evaluate classification performance on the
 aforementioned datasets, we implemented three key base-
 lines using transformer architectures on all datasets: **ViT**
 (Dosovitskiy et al., 2021) without self-supervision; **MAE**
 (He et al., 2022) with masked self-supervision on the in-
 put space; and **I-JEPA** (Assran et al., 2023) with self-
 supervision on the representation space. Beyond that, we

330
 331 **Table 1.** Classification performance on SEED and SEED-IV (Accuracy/F1 Score (%)). SD: Subject Dependent; SI: Subject Independent.
 332 **BOLD:** the best performance are bolded. **UNDERLINE:** the second-best performance.

MODEL	SSL	T	C	R	SEED		SEED-IV	
					SD	SI	SD	SI
DGCNN	×	×	✓	×	75.62/76.19	72.17/71.10	68.89/68.17	67.42/67.74
RGNN	×	×	✓	×	77.26/78.32	73.20/73.12	69.10/69.24	67.34/68.38
TSCEPTION	×	×	✗	✓	74.65/76.02	72.86/74.47	66.83/65.94	69.86/70.08
EEG CONFORMER	×	✓	✗	✗	76.07/76.64	74.75/73.94	71.98/72.46	69.01/68.56
LGGNET	×	✗	✓	✓	78.14/78.92	75.09/76.86	72.14/73.22	<u>72.76/72.06</u>
MMM	✓	✗	✓	✓	77.14/78.61	76.17/77.95	72.24/73.21	71.28/71.80
EPNNE	✗	✗	✗	✗	78.17/77.92	76.47/76.31	72.81/73.46	71.42/71.17
ViT	✗	✓	✓	✗	77.10/78.32	72.94/73.06	71.70/72.83	68.12/67.17
MAE	✓	✓	✓	✗	80.27/81.78	75.73/75.85	73.29/74.48	70.64/70.15
I-JEPA	✓	✓	✓	✗	80.81/81.30	<u>76.93/77.84</u>	74.73/75.04	<u>72.24/71.47</u>
RCT-MNM	✓	✓	✓	✓	82.14/82.10	77.91/78.03	<u>74.46/74.81</u>	72.90/73.44

348 **Table 2.** Fine-tuning classification AUROC (%) on MSIT/ECR.
 349 **BOLD:** the best performance are bolded. **UNDERLINE:** the
 350 second-best performance.

MODEL	SSL	T	C	R	MSIT	ECR
SVM	✗	✗	✗	✗	80.87	80.10
DGCNN	✗	✗	✓	✗	81.72	82.69
RGNN	✗	✗	✓	✗	82.39	81.76
LGG-NET	✗	✗	✓	✓	83.89	84.13
ViT	✗	✓	✓	✗	82.88	83.91
MAE	✓	✓	✓	✗	86.22	86.67
I-JEPA	✓	✓	✓	✗	87.17	<u>86.59</u>
RCT-MNM	✓	✓	✓	✓	89.34	89.05

363 implemented other baselines widely used in emotion recognition: **DGCNN** (Song et al., 2018), **RGNN** (Zhong et al., 2020), **TSception** (Ding et al., 2022), **EEG Conformer** (Song et al., 2022), **MMM** (Yi et al., 2024), **LGGNet** (Ding et al., 2023), **EPNNE** (Zhang et al., 2024), and task engagement classification: **SVM** (Provenza et al., 2019). We also employed **DGCNN**, **RGNN**, and **LGGNet** for task engagement classification, since they are graph neural network based models which rather match the inductive bias of more distributed spatial representations in the MSIT, ECR datasets.

4.2. Experimental Results

377 Table 2 presents the fine-tuning AUROC results for the
 378 MSIT and ECR tasks. Traditional models such as SVM per-
 379 form the worst, lacking the ability to capture complex neural
 380 dynamics. Graph-based models, including DGCNN, RGNN,
 381 and LGG-Net, improve performance by incorporating spa-
 382 tial connectivity but still fall short of transformer-based ap-
 383 proaches. Among transformer-based models, ViT provides
 384

a baseline for self-attention mechanisms, while MAE and I-JEPA leverage self-supervised learning for improved representations. Our RCT-MNM model outperforms all baselines (89.34%, 89.05%), benefiting from spatial-temporal modeling and region-channel self-attention, which enables it to capture fine-grained sEEG representations crucial for cognitive conflict resolution. Overall, our model demonstrates its ability to effectively leverage self-supervised learning (SSL) to handle both subject-dependent and subject-independent classification tasks, achieving consistent improvements across datasets and metrics. The results highlight the robustness of our approach in capturing both intra- and inter-subject variability, making it a promising solution for emotion recognition and related applications. Table 1 presents classification results on SEED and SEED-IV under subject-dependent and subject-independent settings. Our model, RCT-MNM, consistently outperforms baselines. Table 3 further validates our approach on DEAP dataset. Across all evaluations, RCT-MNM surpasses both spatial modeling architecture (DGCNN, RGNN, LGG-Net) and purely temporal models (EEG Conformer) by effectively balancing temporal, regional, and channel interactions. The hierarchical temporal-first encoding ensures stable representations before spatial aggregation, while region-channel self-attention refines functional connectivity patterns. Compared with other self-supervised learning models, RCT-MNM further enhance the pre-training process with RCT-Mask with strong spatial inductive bias and adapt for downstream tasks using Adaptive Refinement process. Additionally, the model's ability to generalize across subjects suggests its effectiveness in handling inter-subject variability, a key challenge in EEG-based cognitive and emotion analysis. Overall, RCT-MNM advances neural representation learning by integrating hierarchical attention, self-supervised learning, and regional conditioning, achieving superior performance in emotion recognition and cognitive task engagement.

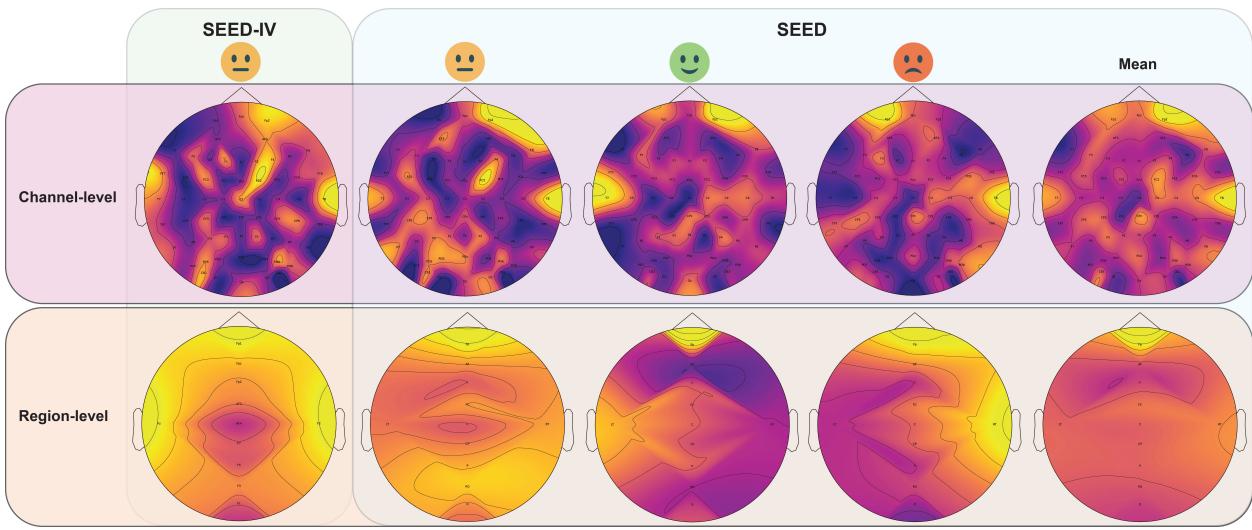


Figure 2. **Region- and Channel-Level Class Activation Maps.** This figure visualizes EEG activation patterns for different emotions in SEED and SEED-IV datasets, showing channel-level fine-grained activations and region-level aggregated representations for a holistic view of emotional processing.

4.3. Interpretations and Visualizations

Figure 2 presents class activation maps (Alarcao & Fonseca, 2017) at the channel and region levels for EEG-based emotion recognition across the SEED (Neutral, Positive, Negative) and SEED-IV (Neutral) datasets. The rightmost column (Mean) shows the average activation across conditions for SEED. The channel-level CAMs reveal distinct spatial distributions of neural activity across emotions. For SEED-IV, the model primarily focuses on frontal and central regions in the neutral condition. For SEED, neutral emotions exhibit similar activations as for SEED-IV. Positive emotions show diffused activation patterns particularly in parietal and occipital regions, and negative emotions elicit heightened activations in the frontal and temporal regions, aligning with findings that these areas are associated with negative affect and stress processing. At the region level, a more structured and spatially smooth representation of neural activity emerges, demonstrating the effectiveness of regional aggregation in our RCT-MNM framework. For SEED-IV (Neutral), activations are dominant in the frontal and central regions, reinforcing their role in neutral emotional states. In SEED, positive emotions show higher activations in the posterior regions, while negative emotions show stronger frontal dominance, consistent with the valence asymmetry hypothesis, which associates left-frontal activity with positive emotions and right-frontal activity with negative emotions. The region-level activations appear smoother and more interpretable compared to channel-level activations, suggesting that hierarchical spatial modeling enhances feature abstraction and generalization.

5. Conclusion

In this paper, we propose Region-Channel-Temporal Masked Neural Modeling (RCT-MNM), a self-supervised learning framework for universal EEG/sEEG representation learning. RCT-MNM leverages region-channel-temporal self-attention and masked neural modeling to effectively capture temporal and spatial dependencies in EEG and sEEG signals. By employing regional conditioned self-supervised learning (RC), the model learns robust neural representations that enhance both emotion recognition and psychiatric disorder analysis. We introduce a hierarchical structure, where local channel-level features are aggregated into region-level embeddings, improving both efficiency and interpretability. Experimental results on multiple EEG datasets for emotion recognition (SEED, SEED-IV, and DEAP) and sEEG datasets for task engagement classification (MSIT and ECR) demonstrate that RCT-MNM outperforms state-of-the-art models in both subject-dependent and subject-independent settings. The visualization of class activation maps and self-attention graphs further validates the model’s ability to learn meaningful neurophysiological representations. Future work will explore real-time BCI applications, cross-domain transfer learning, and expanded pre-training datasets to further enhance model scalability and generalization.

Impact Statement

This work explores self-supervised learning for EEG/sEEG representation by integrating region-channel-temporal self-

attention and masked neural modeling to enhance feature extraction for emotion recognition and psychiatric analysis. Our approach improves the generalizability and interpretability of EEG/sEEG-based models, with potential applications in mental health monitoring and brain-computer interfaces (BCIs). There are no immediate negative societal impacts anticipated.

References

- Alarcao, S. M. and Fonseca, M. J. Emotions recognition using eeg signals: A survey. *IEEE transactions on affective computing*, 10(3):374–393, 2017.
- Assran, M., Duval, Q., Misra, I., Bojanowski, P., Vincent, P., Rabbat, M., LeCun, Y., and Ballas, N. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15619–15629, 2023.
- Bardes, A., Ponce, J., and Lecun, Y. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. In *ICLR 2022-International Conference on Learning Representations*, 2022.
- Chien, H.-Y. S., Goh, H., Sandino, C. M., and Cheng, J. Y. Maeeg: Masked auto-encoder for eeg representation learning. In *NeurIPS Workshop*, 2022. URL <https://arxiv.org/abs/2211.02625>.
- Ding, Y., Robinson, N., Zhang, S., Zeng, Q., and Guan, C. Tsception: Capturing temporal dynamics and spatial asymmetry from eeg for emotion recognition. *IEEE Transactions on Affective Computing*, 14(3):2238–2250, 2022.
- Ding, Y., Robinson, N., Tong, C., Zeng, Q., and Guan, C. Lggnnet: Learning from local-global-graph representations for brain–computer interface. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Drane, D. L., Pedersen, N. P., Sabsevitz, D. S., Block, C., Dickey, A. S., Alwaki, A., and Kheder, A. Cognitive and emotional mapping with seeg. *Frontiers in Neurology*, 12:627981, 2021.
- Gao, Z., Wang, X., Yang, Y., Mu, C., Cai, Q., Dang, W., and Zuo, S. Eeg-based spatio-temporal convolutional neural network for driver fatigue evaluation. *IEEE transactions on neural networks and learning systems*, 30(9):2755–2763, 2019.
- Gevins, A., Leong, H., Smith, M. E., Le, J., and Du, R. Mapping cognitive brain function with modern high-resolution electroencephalography. *Trends in Neurosciences*, 18(10):429–436, 1995. ISSN 0166-2236. doi: [https://doi.org/10.1016/0166-2236\(95\)94489-R](https://doi.org/10.1016/0166-2236(95)94489-R). URL <https://www.sciencedirect.com/science/article/pii/016622369594489R>.
- Goschke, T. Dysfunctions of decision-making and cognitive control as transdiagnostic mechanisms of mental disorders: Advances, gaps, and needs in current research. *International journal of methods in psychiatric research*, 23(S1):41–57, 2014.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Houssein, E. H., Hammad, A., and Ali, A. A. Human emotion recognition from eeg-based brain–computer interface using machine learning: a comprehensive review. *Neural Computing and Applications*, 34(15):12527–12557, 2022.
- Jiang, W., Zhao, L., and Lu, B.-l. Large brain model for learning generic representations with tremendous eeg data in bci. In *The Twelfth International Conference on Learning Representations*, 2024.
- Koelstra, S., Muhl, C., Soleymani, M., Lee, J.-S., Yazdani, A., Ebrahimi, T., Pun, T., Nijholt, A., and Patras, I. Deep: A database for emotion analysis; using physiological signals. *IEEE transactions on affective computing*, 3(1):18–31, 2011.
- Kret, M. E. and Ploeger, A. Emotion processing deficits: a liability spectrum providing insight into comorbidity of mental disorders. *Neuroscience & Biobehavioral Reviews*, 52:153–171, 2015.
- Lachaux, J., Rudrauf, D., and Kahane, P. Intracranial eeg and human brain mapping. *Journal of Physiology-Paris*, 97(4):613–628, 2003. ISSN 0928-4257. doi: <https://doi.org/10.1016/j.jphysparis.2004.01.018>. URL <https://www.sciencedirect.com/science/article/pii/S0928425704000348>.
- Neuroscience and Computation.
- Liu, J., Younk, R., Drahos, L. M., Nagrale, S. S., Yadav, S., Widge, A. S., and Shoaran, M. Neural decoding and feature selection methods for closed-loop control of avoidance behavior. *Journal of Neural Engineering*, 21(5):056041, 2024.

- 495 Mentzelopoulos, G., Chatzipantazis, E., Ramayya, A. G.,
 496 Hedlund, M., Buch, V., Daniilidis, K., Kording, K., and
 497 Vitale, F. Neural decoding from stereotactic eeg: ac-
 498 counting for electrode variability across subjects. In *The*
 499 *Thirty-eighth Annual Conference on Neural Information*
 500 *Processing Systems*, 2024.
- 501 Ng, M. Y. and Weisz, J. R. Annual research review: Building
 502 a science of personalized intervention for youth mental
 503 health. *Journal of Child Psychology and Psychiatry*, 57
 504 (3):216–236, 2016.
- 505 Organization, W. H. The world health report 2001: Mental
 506 health: new understanding, new hope. 2001.
- 507 Peled, N., Gholipour, T., Paultk, A., Felsenstein, O.,
 508 Dougherty, D., Widge, A., Eskandar, E., Cash, S.,
 509 Hamalainen, M., and Stufflebeam, S. Invasive electrodes
 510 identification and labeling. *GitHub Repos*, 10, 2017.
- 511 Provenza, N. R., Paultk, A. C., Peled, N., Restrepo, M. I.,
 512 Cash, S. S., Dougherty, D. D., Eskandar, E. N., Borton,
 513 D. A., and Widge, A. S. Decoding task engagement
 514 from distributed network electrophysiology in humans.
 515 *Journal of neural engineering*, 16(5):056015, 2019.
- 516 Rafiei, M. H., Gauthier, L. V., Adeli, H., and Takabi, D. Self-
 517 supervised learning for electroencephalography. *IEEE*
 518 *Transactions on Neural Networks and Learning Systems*,
 519 35(2):1457–1471, 2024. doi: 10.1109/TNNLS.2022.
 520 3190448.
- 521 Rouast, P. V., Adam, M. T., and Chiong, R. Deep learning
 522 for human affect recognition: Insights and new develop-
 523 ments. *IEEE Transactions on Affective Computing*, 12
 524 (2):524–543, 2019.
- 525 Saha, S. and Baumert, M. Intra- and inter-subject variability
 526 in eeg-based sensorimotor brain computer interface: A re-
 527 view. *Frontiers in Computational Neuroscience*, 13, 2020.
 528 ISSN 1662-5188. doi: 10.3389/fncom.2019.00087. URL
 529 <https://www.frontiersin.org/journals/computational-neuroscience/articles/10.3389/fncom.2019.00087>.
- 530 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R.,
 531 Parikh, D., and Batra, D. Grad-cam: Visual explana-
 532 tions from deep networks via gradient-based localiza-
 533 tion. In *Proceedings of the IEEE international conference on*
 534 *computer vision*, pp. 618–626, 2017.
- 535 Shazeer, N. GLU variants improve transformer. *CoRR*,
 536 abs/2002.05202, 2020. URL <https://arxiv.org/abs/2002.05202>.
- 537 Song, T., Zheng, W., Song, P., and Cui, Z. Eeg emotion
 538 recognition using dynamical graph convolutional neural
 539 networks. *IEEE Transactions on Affective Computing*, 11
 540 (3):532–541, 2018.
- 541 Song, Y., Zheng, Q., Liu, B., and Gao, X. Eeg conformer:
 542 Convolutional transformer for eeg decoding and visu-
 543 alization. *IEEE Transactions on Neural Systems and*
 544 *Rehabilitation Engineering*, 31:710–719, 2022.
- 545 Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., and Liu, Y. Roformer: Enhanced transformer with rotary position
 546 embedding. *Neurocomputing*, 568:127063, 2024.
- 547 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
 548 L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. Atten-
 549 tion is all you need. In Guyon, I., Luxburg, U. V.,
 550 Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S.,
 551 and Garnett, R. (eds.), *Advances in Neural Information*
 552 *Processing Systems*, volume 30. Curran Associates, Inc.,
 553 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fb053c1c4a845aa-Paper.pdf.
- 554 Wang, G., Liu, W., He, Y., Xu, C., Ma, L., and Li, H. Eegpt: Pretrained transformer for universal and reliable
 555 representation of eeg signals. In *The Thirty-eighth Annual*
 556 *Conference on Neural Information Processing Systems*, 2024.
- 557 Ye, M., Chen, C. P., and Zhang, T. Hierarchical dynamic
 558 graph convolutional network with interpretability for eeg-
 559 based emotion recognition. *IEEE transactions on neural*
 560 *networks and learning systems*, 2022.
- 561 Yi, K., Wang, Y., Ren, K., and Li, D. Learning topology-
 562 agnostic eeg representations with geometry-aware model-
 563 ing. *Advances in Neural Information Processing Systems*,
 564 36, 2024.
- 565 Zhang, H., Zuo, T., Chen, Z., Wang, X., and Sun, P. Z. Evolutionary ensemble learning for eeg-based cross-
 566 subject emotion recognition. *IEEE Journal of Biomedical*
 567 *and Health Informatics*, 28(7):3872–3881, 2024. doi:
 568 10.1109/JBHI.2024.3384816.
- 569 Zheng, W.-L. and Lu, B.-L. Investigating critical frequency
 570 bands and channels for eeg-based emotion recognition
 571 with deep neural networks. *IEEE Transactions on au-
 572 tonomous mental development*, 7(3):162–175, 2015.
- 573 Zheng, W.-L., Liu, W., Lu, Y., Lu, B.-L., and Cichocki, A. Emotionmeter: A multimodal framework for recognizing
 574 human emotions. *IEEE transactions on cybernetics*, 49
 575 (3):1110–1122, 2018.
- 576 Zhong, P., Wang, D., and Miao, C. Eeg-based emotion
 577 recognition using regularized graph neural networks.
 578 *IEEE Transactions on Affective Computing*, 13(3):1290–
 579 1301, 2020.

A. Additional Classification Results

Table 3. Classification performance on DEAP (Valence/Arousal) (%). **BOLD**: the best performance are bolded. UNDERLINE: the second-best performance.

MODEL	SUBJECT-DEPENDENT		SUBJECT-INDEPENDENT	
	ACCURACY	F1 SCORE	ACCURACY	F1 SCORE
DGCNN	64.97/58.52	68.08/60.13	61.76/64.82	60.15/62.92
RGNN	65.65/67.54	64.70/58.74	63.45/65.26	62.28/66.11
TSCEPTION	64.19/66.26	63.86/65.47	62.83/64.78	63.86/ 68.08
EEG CONFORMER	62.48/65.02	64.09/62.60	62.00/60.16	61.76/57.84
LGGNET	67.28/66.61	67.14/66.69	66.17/65.85	64.54/67.44
MMM	65.18/67.01	64.96/ <u>68.97</u>	66.03/66.16	65.58/65.64
EPNNE	67.20/64.23	65.06/66.17	65.58/63.31	65.00/63.88
ViT	66.40/63.27	65.83/64.42	65.15/64.67	64.67/63.52
MAE	67.20/68.17	66.01/67.16	65.58/64.52	65.30/65.27
I-JEPA	<u>69.22/69.01</u>	<u>68.19/68.29</u>	<u>66.70/67.52</u>	<u>67.56/66.47</u>
RCT-MNM	71.27/70.92	71.78/70.60	67.47/67.77	68.15/67.94

B. Dataset Details

MSIT/ECR. The MSIT and ECR datasets are utilized to study human cognitive and emotional conflict responses. These datasets consist of electrophysiological recordings from 17 participants with pharmaco-resistant complex partial seizures undergoing invasive monitoring. Participants performed two distinct conflict-based tasks: the Multi-Source Interference Task (MSIT) and the Emotional Conflict Resolution (ECR) task. The MSIT involved identifying a target number among distractors under varying levels of congruence, while the ECR task required resolving emotional conflict between facial expressions and overlaid text. Recordings were made using depth electrodes placed in up to 30 brain regions, capturing local field potentials (LFPs) at 2 kHz and then downsampled to 1 kHz. The number of electrodes used for recording varied across participants from 37 to 195. The LFP recordings from different channels were localized to standard brain regions using electrode mapping tools. To remove the low and high-frequency noise, a band-pass filter from 1.0-150Hz was applied to the original LFP. Line noise was removed using notch filters on 60 Hz and its harmonics. A data augmentation step by splitting each trial into smaller non-overlapping 4-second segments was applied. Then, the neural data $\mathbf{X} \in \mathbb{R}^{C \times 4000}$ was used to predict rest-state versus task-state as $y \in \{0, 1\}$, where C is the number of channels. We split each session into the training-validation-testing sets in a chronological manner with a 80%-10%-10% ratio, between them we left 4-second segments to make sure no adjacent segments will appear in both training and testing data. This dataset provides insights into task-specific and generalized cognitive processes by leveraging robust, multi-site neural recordings, with applications in psychiatric adaptive deep brain stimulation and beyond.

DEAP. The DEAP dataset is designed for analyzing human affective states using multimodal data. It includes EEG and peripheral physiological signals from 32 participants as they watched 40 one-minute music video clips. Participants rated the videos on arousal, valence, dominance, liking, and familiarity. EEG data was recorded using 32 electrodes at 512 Hz, downsampled to 128 Hz. The 3 seconds pre-trial baseline was removed for each trial. To remove the low and high-frequency noise, a band-pass filter from 4.0-45Hz was applied to the original EEG. The class label for each dimension is from 1 to 9, hence 5 was selected as a threshold to project the 9 discrete values into low and high classes in each dimension. A data augmentation step by splitting each trial into smaller non-overlapping 4-second segments was applied. Then, the neural data $\mathbf{X} \in \mathbb{R}^{32 \times 512}$ was used to predict low valence/arousal versus high ones as $y \in \{0, 1\}$. There are two types of experiment settings in this paper: I) subject-dependent: trial-wise 10-fold cross-validation; II) subject-independent: leave-one-subject-out cross-validation. Randomly shuffling the segments among trials before the training-testing split of the data could make the adjacent segments be in training and testing data, which will give high classification results. But the accuracy will drop when the highly correlated segments are never seen by the model in the real-world situation. To get the more generalized evaluation, the folds are split among trials, which will make sure the adjacent segments in one trial will not appear in both training and testing data. The ratio of training to validation data is 80%-20%.

SEED. The SEED dataset is a discrete EEG emotion dataset collected through video stimuli. It includes EEG recordings from fifteen subjects, captured via 62 sensors while participants viewed multiple Chinese film clips, each eliciting one of

605 three emotional responses, namely positive, neutral, and negative. EEG data was recorded at 1 kHz, downsampled to 200 Hz.
 606 For each subject, recordings were conducted across three separate sessions, with each session containing EEG data from
 607 fifteen film clips (fifteen trials). Additionally, the SEED-IV dataset, an extended version of SEED, comprises EEG data
 608 from the same fifteen subjects using the same 62-channel setup. Each subject's data in SEED-IV spans three sessions, with
 609 each session including 24 trials, and each trial elicited one of four emotion responses, namely happy, neutral, sad, and fear.
 610 A data augmentation step by splitting each trial into smaller non-overlapping 4-second segments was applied. Then, the
 611 neural data $\mathbf{X} \in \mathbb{R}^{62 \times 800}$ was used to predict distinct emotion responses as $y \in \{0, 1, 2\}$ for SEED and $y \in \{0, 1, 2, 3\}$ for
 612 SEED-IV. There are two types of experiment I) subject-dependent: trial-wise 5-fold cross-validation; II) subject-independent:
 613 leave-one-subject-out cross-validation as DEAP. The ratio of training to validation data is 80%-20%.

615 C. Details of Baseline Models

616 **DGCNN** (Song et al., 2018) is a deep learning framework for EEG emotion recognition that utilizes dynamical graph
 617 convolutional neural networks. It dynamically constructs an adjacency matrix to capture the intrinsic relationships between
 618 EEG channels, enhancing feature extraction. The model applies graph convolution to extract localized spatial dependencies,
 619 followed by a 1×1 convolution layer for feature transformation. A fully connected layer is used for classification. Unlike
 620 traditional GCNNs, DGCNN learns the adjacency matrix adaptively during training, improving EEG emotion recognition
 621 performance.

622 **RGNN** (Zhong et al., 2020) is a deep learning model designed for EEG-based emotion recognition by incorporating
 623 biologically inspired graph structures. It models inter-channel relationships via an adjacency matrix, leveraging neuroscience
 624 principles to capture both local and global spatial dependencies. RGNN introduces two key regularizers: Node-wise
 625 Domain Adversarial Training (NodeDAT) to enhance cross-subject generalization and Emotion-Aware Distribution Learning
 626 (EmotionDL) to mitigate label noise. The model extends simple graph convolution networks (SGC) for efficient EEG
 627 representation learning.

628 **TSeption** (Ding et al., 2022) is a multi-scale convolutional neural network designed for EEG-based emotion recognition.
 629 It consists of three key layers: a dynamic temporal layer, which employs multi-scale 1D convolutional kernels to capture
 630 temporal dynamics and frequency representations; an asymmetric spatial layer, which leverages the brain's hemispheric
 631 asymmetry to extract global and hemisphere-specific spatial patterns; and a high-level fusion layer, which integrates
 632 learned spatial and temporal representations. TSeption effectively captures EEG temporal and spatial features, improving
 633 classification performance in emotion recognition tasks while maintaining a compact and efficient model architecture.

634 **EEG Conformer** (Song et al., 2022) is a hybrid convolutional-transformer model designed for EEG decoding and
 635 visualization. It consists of three key components: a convolution module for capturing low-level local temporal and
 636 spatial features, a self-attention module for modeling global temporal dependencies, and a classifier module based on
 637 fully connected layers for final prediction. The convolutional layers extract localized features using temporal and spatial
 638 convolutions, while an average pooling layer reduces redundant information. The self-attention mechanism is applied to the
 639 transformed feature maps to learn long-term dependencies.

640 **MMM** (Yi et al., 2024) is a self-supervised pre-training framework designed for EEG-based emotion recognition. It unifies
 641 EEG channel topologies using a Multi-dimensional position encoding, Multi-level channel hierarchy, and a Multi-stage
 642 pre-training strategy to obtain topology-agnostic representations. The model employs a masked autoencoder (MAE)-based
 643 approach, where EEG signals are mapped onto a unified topology that ensures cross-dataset generalization. Region-wise
 644 tokens are introduced to capture local spatial information, while a structured masking strategy (random and region-wise
 645 masking) enhances representation robustness.

646 **LGG-Net** (Ding et al., 2023) is a neurologically inspired graph neural network (GNN) designed for brain-computer interface
 647 (BCI) applications using EEG data. It integrates local-global-graph (LGG) representations to capture both localized and
 648 global brain activities. The model consists of a temporal learning block employing multi-scale 1D convolutions for dynamic
 649 EEG feature extraction and a graph learning block with local- and global-graph-filtering layers to model brain region
 650 interactions. LGG-Net applies neurophysiological priors, defining EEG as a structured LGG with functional subgraphs, and
 651 learns hierarchical spatial dependencies through an instance-specific global adjacency matrix.

652 **EPNNE** (Zhang et al., 2024) is an end-to-end framework designed for cross-subject EEG-based emotion recognition. It
 653 employs an evolutionary programming (EP)-based optimization strategy to automatically search for an optimal neural
 654 network architecture, mitigating inter-subject variability in EEG signals. The model integrates an ensemble of multilayer
 655

perceptrons (MLPs) optimized through crossover and mutation operations, enhancing the generalizability of emotion recognition.

D. Metrics

We evaluated the performance of all baseline models and our model using the Area Under the Receiver Operating Characteristic Curve (AUROC) for task engagement classification on MSIT and ECR, and using Accuracy and Weighted F1-Score for emotion recognition on SEED, SEED-IV, and DEAP.

E. Additional Interpretations and Visualizations

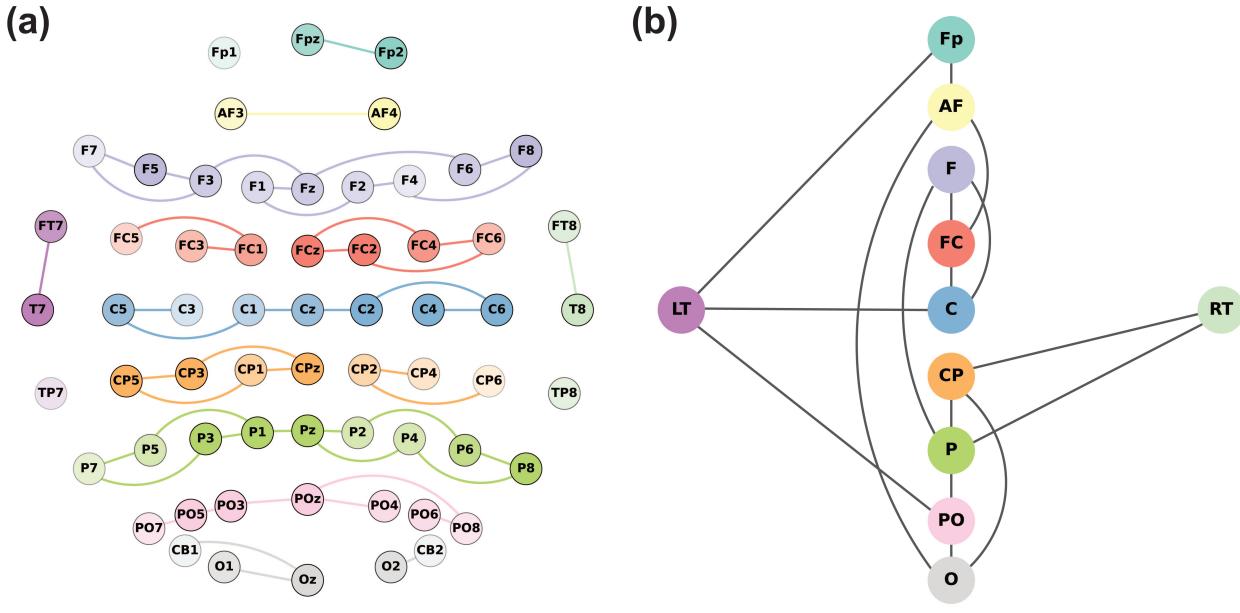


Figure 3. Region-Channel Self-Attention Graphs. (a) Localized EEG channel attention within regions. (b) Region-level attention dependencies, with nodes as aggregated regional embeddings. Connections highlight strong learned attention relationships between channels (colored) and regions (gray-scaled). The color intensity of channels in (a) indicates its level of attention to its corresponding region.

Figure 3 illustrates the Region-Channel Self-Attention Graphs for SEED dataset, providing insight into the structured spatial dependencies learned by the Region-Channel Transformer with Masked Neural Modeling framework. The channel connectivity graph in Figure 3(a) depicts the learned intra-regional, functional connections between channels measured using attention scores for local feature extraction. In Figure 3(b), each node of the channel connectivity graph represents a major brain region, with edges denoting the learned attention relationships between different regions. It is observed that there is stronger connectivity among adjacent and functionally related regions, such as frontal-central (FC-C) and parietal-occipital (P-PO) connections. Long-range dependencies between temporal (LT, RT) and pre-frontal (Fp) or parietal (P, PO) regions, highlight interactions crucial for cognitive and affective processes. Asymmetric attention distribution is observed which aligns with known emotional and neuropsychiatric functions.

Overall, the visualization results validate the RCT-MNM framework's ability to capture both local and global neural dependencies, leading to robust and interpretable emotion classification. The hierarchical self-attention mechanism effectively integrates channel- and region-level features, facilitating better generalization and alignment with functional brain networks involved in emotional and cognitive processes.

F. Additional Ablation Study

Table 4. Ablations to validate the effectiveness of MNM-Loss on MSIT/ECR (AUROC (%)) and SEED (Accuracy/F1 Score (%)).

$\mathcal{L}_{\text{input}}$	\mathcal{L}_{rep}	\mathcal{L}_{reg}	MSIT	ECR	SEED	
					DEPENDENT	INDEPENDENT
✓	✓	✓	89.34	89.05	82.14/82.10	77.91/78.03
✓	✓	✗	89.08	87.70	81.88/82.03	76.86/77.93
✓	✗	✗	87.26	87.02	81.74/81.44	75.66/76.93
✗	✓	✓	88.23	88.64	81.98/ 82.57	77.76/77.84

Table 5. Ablations to validate the effectiveness of Region-Channeel-Temporal Masking (RCT-Mask), Regioal Conditioning (RC), and Adaptive Attention-Based Refining on MSIT/ECR (AUROC (%)) and SEED (Accuracy/F1 Score (%)).

MODEL	MSIT	ECR	SEED	
			DEPENDENT	INDEPENDENT
OUR	89.34	89.05	82.14/82.10	77.91/78.03
w/o RC	88.37	87.77	81.85/ 82.22	77.24/77.21
RCT-MASK → RANDOM	88.24	87.12	81.05/81.38	77.09/76.14
w/o REFINE	87.71	87.20	81.47/81.11	76.64/76.81

We performed ablation studies to evaluate the effectiveness of different components in our proposed framework. Specifically, we examined the contributions of the loss components in the RCT-MNM in Table 4 and key architectural components, including Region-Channel-Temporal masking (RCT-Mask), Regional Conditioning (RC), and Adaptive Attention-Based Layer Refining in Table 5. The results highlight the importance of these design choices in achieving state-of-the-art performance across multiple datasets, including MSIT, ECR, and SEED.

Loss Component. We systematically remove one or more loss components. The complete model, utilizing all three losses, achieves the relatively higher performance across datasets. Removing the VICReg loss leads to a notable performance drop across all datasets except MSIT. It enforces a balanced representation space by penalizing covariance redundancy and maximizing feature invariance. This contributes to better generalization, especially for datasets with high inter-subject variability, such as SEED for subject-independent emotion recognition. When the embedding loss is excluded, the performance degradation is more pronounced. The representatioin alignment loss ensures consistency between masked embeddings and their corresponding latent representations, which is crucial for retaining task-relevant information. Its absence makes the model less robust, especially for SEED. Removing input reconstruction loss reduces Accuracy/F1 Score for SEED. This indicates that input loss aligns reconstructed inputs with the original signals, contributing to stability but is partially redundant when the embedding and VICReg losses are present.

Masking. Replacing the RCT-Mask with random masking significantly reduces performance. The RCT-Mask ensures structured masking along region, channel and temporal dimensions, which encourages the model to learn meaningful relationships across these domains. Random masking disrupts this structure, leading to degraded representations and poor generalization.

Regional Conditioning. Removing RC results in performance degradation. RC enables task-specific contextualization of target embeddings, which improves alignment with downstream objectives. While the drop is less severe than other ablations, it demonstrates the importance of context-aware feature learning.

Refinement. Removing the refinement step leads to consistent performance drops across all datasets. Refinement fine-tunes the intermediate-layer encoder representations to align with downstream tasks, significantly boosting performance. Its absence results in less effective task adaptation.

The results demonstrate the synergistic contributions of our proposed components. The **RCT-Mask** ensures the model learns meaningful spatial-temporal relationships, while **RC** improves feature contextualization. **Loss regularization**, particularly \mathcal{L}_{reg} , enhances the representation space's robustness, and **Refinement** aligns these representations with specific tasks. Collectively, these components enable our framework to achieve state-of-the-art performance on MSIT, ECR, and SEED

770 datasets, demonstrating its effectiveness for emotion recognition and psychiatric disorder prediction.
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824

G. Additional Figure

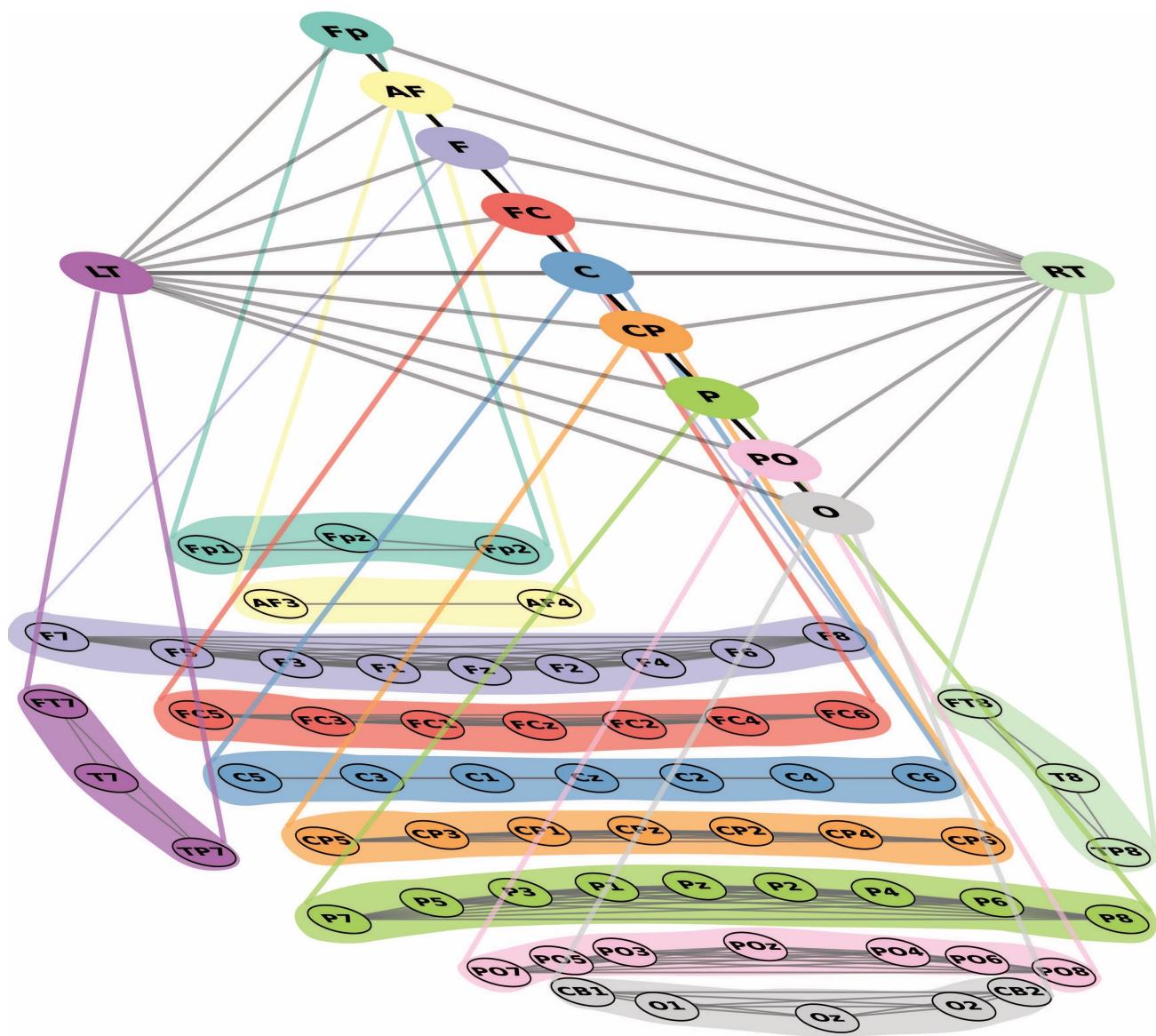


Figure 4. Region-Channel Self-Attention Map.