

# EEG\_RL-Net: Enhancing EEG MI Classification through Reinforcement Learning-Optimised Graph Neural Networks

Htoo Wai Aung, Jiao Jiao Li, Yang An, and Steven W. Su\*, *Senior Member, IEEE*

**Abstract**—Brain-Computer Interfaces (BCIs) rely on accurately decoding electroencephalography (EEG) motor imagery (MI) signals for effective device control. Graph Neural Networks (GNNs) outperform Convolutional Neural Networks (CNNs) in this regard, by leveraging the spatial relationships between EEG electrodes through adjacency matrices. The EEG\_GLT-Net framework, featuring the state-of-the-art EEG\_GLT adjacency matrix method, has notably enhanced EEG MI signal classification, evidenced by an average accuracy of 83.95% across 20 subjects on the PhysioNet dataset. This significantly exceeds the 76.10% accuracy rate achieved using the Pearson Correlation Coefficient (PCC) method within the same framework.

In this research, we advance the field by applying a Reinforcement Learning (RL) approach to the classification of EEG MI signals. Our innovative method empowers the RL agent, enabling not only the classification of EEG MI data points with higher accuracy, but effective identification of EEG MI data points that are less distinct. We present the EEG\_RL-Net, an enhancement of the EEG\_GLT-Net framework, which incorporates the trained EEG\_GCN Block from EEG\_GLT-Net at an adjacency matrix density of 13.39% alongside the RL-centric Dueling Deep Q Network (Dueling DQN) block. The EEG\_RL-Net model showcases exceptional classification performance, achieving an unprecedented average accuracy of 96.40% across 20 subjects within 25 milliseconds. This model illustrates the transformative effect of the RL in EEG MI time point classification.

**Index Terms**—Brain-Computer Interfaces (BCIs), Electroencephalography Motor Imagery (EEG MI), Spectral Graph Convolutional Neural Networks (GCNs), Reinforcement Learning (RL), Dueling Deep Q Network (Dueling DQN)

## I. INTRODUCTION

**B**RAIN-COMPUTER INTERFACES establish a connection between the brain and external control devices. Originally developed to assist individuals with motor impairments [1], BCIs translate brain signals acquired through measurements such as electrocorticography (ECoG) and electroencephalogram (EEG) into actionable commands for electronic control devices including wheelchairs and exoskeleton robots.

Htoo Wai Aung is the first author of this paper, and he is with the School of Biomedical Engineering, Faculty of Engineering and IT, University of Technology Sydney, NSW 2007, Australia (e-mail: htoowai.aung@student.uts.edu.au).

Jiao Jiao Li is with the School of Biomedical Engineering, Faculty of Engineering and IT, University of Technology Sydney, NSW 2007, Australia (e-mail: jiaojiao.li@uts.edu.au).

Yang An is with the Faculty of Engineering and IT, University of Technology Sydney, NSW 2007, Australia (e-mail: yang.an-1@student.uts.edu.au).

Steven W. Su is with the Faculty of Engineering and IT, University of Technology Sydney, NSW 2007, Australia (e-mail: steven.su@uts.edu.au).

\*Corresponding author.

Although ECoG offers superior signal quality over EEG, its application in BCIs is limited due to invasive route of acquisition, requiring the placement of electrodes directly on the cerebral cortex [2]. Meanwhile, EEG is a much more accessible and hence popular signal acquisition method as it involves non-invasive placement of electrodes on the scalp. EEG is widely used to record various types of brain signals, from spontaneous and stimulus-evoked signals to event-related potentials [3]. Its clinically relevant applications extend to dementia classification [4], depression state assessment [5], seizure detection [6], and the classification of cognitive and motor tasks [7], including motor imagery (MI) tasks [5], [8], [9].

MI involves the mental simulation of motor actions, such as movements of the hands, feet, or tongue, without performing the physical movements [10], [11]. This technique is crucial in neuroscience and rehabilitation, with real-world relevance especially for individuals with motor impairments, such as stroke survivors. Through integration with an external control device, MI enables the physically impaired to perform daily activities that are not otherwise possible, leading to potentially life-changing benefits by improving quality of life and reducing the level of chronic care. By integrating MI and BCIs, EEG based MI signals can be decoded and used to control external devices, enabling real-time feedback and facilitating patient-intended movements through accurate signal interpretation [12].

Deep learning, a subset of machine learning, utilises multiple layers of neural networks to process a variety of data forms. Convolutional Neural Networks (CNNs), which mimic natural image recognition in the human visual system, are part of the deep learning family and excel in computer vision tasks [13]–[15]. However, their application is restricted to Euclidean data, such as 1-dimensional sequences and 2-dimensional grids [15]. CNNs struggle with non-Euclidean data, failing to accurately capture the intrinsic structure and connectivity of the data.

Graph Convolutional Networks (GCNs) have been developed to perform convolutional operations on graphs, which can handle non-Euclidean data due to incorporating topological relationships during convolution. GCNs can represent complex structures and variations in these structures, which may be heterogeneous or homogeneous, weighted or unweighted, signed or unsigned [16]. They support various types of graph analyses, including node-level, edge-level, and graph-level tasks [16], [17]. GCNs are particularly effective at classifying EEG signals as a graph-level task [8], [18]. For this application,

EEG signal readings from each channel are treated as node attributes, and the relationships between EEG electrodes are represented by an adjacency matrix, hence surpassing the capabilities of traditional CNNs.

There are two primary categories of GCNs: spatial [19]–[22] and spectral methods [23]–[25]. Some challenges are encountered with the spatial method [26], [27] especially in matching local neighbourhoods. Both time domain and frequency domain features can be extracted from EEG signals to perform GCN operations [4], [28]–[30]. Frequency domain features include Power Spectral Density (PSD) and Power Ratio (PR) for various bands, such as  $\delta$  (0.5–4Hz),  $\theta$  (4–8Hz),  $\alpha$  (8–13Hz),  $\beta$  (13–30Hz), and  $\gamma$  (30–110Hz) within specified time windows. Time domain features, such as Root Mean Square (RMS), skewness, minmax, variance, number of zero crosses, Hurst Exponent, Petrosian fractal, and Higuchi, are also extracted for GCN operations during specific time windows. These features are integral to window-based GCN methods.

In the GCNs-Net [8], individual time point signals at each channel are treated as distinct features. This method is designed for real-time EEG MI signal classification, focusing on  $\frac{1}{160}s$  time point signals. The constructing of an effective adjacency matrix is crucial for GCN operations, and different methods have been explored in various studies, including: Geodesic method, which relies on geodesic distances between EEG channels [9], [31]–[33]; using Pearson Coefficient Correlation (PCC) to evaluate interchannel correlations [8], [27], [29], [34], [35]; and experimenting with a trainable matrix approach [27], [36].

In the EEG\_GLT-Net [18], a sophisticated algorithm known as the EEG Graph Lottery Ticket (EEG\_GLT) is used to optimise the adjacency matrix by exploring various density levels, inspired from the unified GNN sparsification technique (UGS) [37]. This method represents the current state-of-the-art in adjacency matrix construction, significantly enhancing accuracy, F1 score, and computational efficiency on the EEG MI PhysioNet dataset [38] compared to the PCC and Geodesic methods. However, despite the overall superiority of this method, it remains challenging to classify the EEG MI time points remains challenging for some subjects due to signal ambiguity among different MI tasks at specific time points. Consequently, supervised learning on these subjects involves training that forces classification of all time points.

Reinforcement Learning (RL), another subset of machine learning, enables an RL agent to learn sequential decision-making in dynamic environments to maximise cumulative rewards [39]. RL has been primarily applied in robotics and autonomous systems, which require complex sequential decision-making. Deep RL principles have been applied to optimise feature selection for the Classification with Costly Features (CwCF) problem [40], across various public UCI datasets [41] including miniboone, forest, cifar, wine, and mnist. Others [42] have trained an RL agent to minimise feature extraction costs in classifying electromyography (EMG) signals from UCI datasets [41], although this reduction in features compromised accuracy.

In this paper, we introduce EEG\_RL-Net as a new algo-

rithm, with more advanced capability than existing methods for classifying EEG MI time point signals by combining GNNs and RL. Initially, optimal graph features of EEG MI time point signals are extracted using the best weights and adjacency matrix from an EEG\_GCN block, refined to 13.39% density using the EEG\_GLT algorithm. Subsequently, the RL agent makes sequential decisions within an episode of pre-defined horizon length to accurately classify the EEG MI signals. The main contributions of this study are:

- **EEG\_RL-Net:** A new approach for classifying EEG MI time point signals, using a trained RL agent that determines whether to classify or skip each time point based on GNN features. This method greatly enhances performance accuracy by achieving classification as swiftly as possible within predefined episode lengths.
- **Optimal Reward and Max Episode Length Setting:** We evaluated the accuracy and classification speed under various reward settings and maximum episode lengths for each subject, identifying the optimal combinations for simultaneously achieving high accuracy and efficiency.
- **Performance Validation:** We evaluated the performance of each subject under optimal settings against the state-of-the-art EEG\_GLT-Net with  $m_{g\_GLT}$  matrix and PCC adjacency matrix. Our results showed significant enhancement of accuracy and efficiency on the PhysioNet dataset.

## II. METHODOLOGY

### A. Overview

This project is divided into two distinct parts. The first phase focused on training the EEG\_GLT-Net model, as illustrated in Figure 1, to identify the optimal adjacency matrices and spectral GNN weights across different adjacency matrix density levels, employing Algorithm 1. This phase of training spanned from  $t = 1s$  to  $t = 3s$ . Subsequently, the optimal adjacency matrix and spectral GNN weights, determined at the minimal adjacency matrix density level of 13.39%, were selected for the purpose of extracting graph features.

In the project's second phase, the Multilayer Perceptron (MLP) block within the EEG\_GLT-Net was removed, and in its place, the RL (Reinforcement Learning) block was integrated, resulting in the formation of the EEG\_RL-Net, as depicted in Figure 2. The pre-trained optimal weights of the EEG\_GCN block, such as adjacency matrix and spectral GNN, determined at the lowest adjacency matrix density of 13.39%, were then transferred to the EEG\_GCN component of the EEG\_RL-Net architecture. During this phase, all time points from  $t = 0s$  to  $t = 4s$  were utilised, with these points organised into groups spanning a horizon of 20 states, where each point represented a single state. Within each episode's horizon, the RL agent performed action at every state, based on the graph features generated by the GNN segment. These actions involved classifying the state as belonging to Task 1 through to Task 4, or skipping to the next state (Task 0) if the agent determined that it was not yet prepared to classify.

## B. Dataset Description and Pre-processing

Following the approach of papers [8] and [18], this study employed the PhysioNet EEG MI dataset [38], which comprises EEG recordings from 109 subjects acquired using the international 10-10 system with 64 EEG channels. The dataset is structured around four distinct EEG MI tasks, which involve the subject imagining the actions of:

- Task 1: Opening and closing the left fist.
- Task 2: Opening and closing the right fist.
- Task 3: Opening and closing both fists simultaneously.
- Task 4: Opening and closing both feet.

Each participant completed 84 trials, divided into 3 runs with 7 trials per run for each task type. The duration of each trial's recording was 4 seconds, sampled at 160Hz. In our study analyses were specifically conducted on a subset of 20 subjects, labelled  $S_1$  to  $S_{20}$ . Initially, the raw signals were processed solely through a notch filter at the 50Hz power line frequency to eliminate electrical interference, deliberately avoiding other common filtering or denoising techniques to preserve data integrity. Signals from all 64 channels were utilised, with each channel treated as a node and the signal at each time point considered as the node's feature. Additionally, the signals at each channel were normalised to achieve a mean ( $\mu$ ) of 0 and a standard deviation ( $\sigma$ ) of 1.

## C. Graph Feature Extraction

1) *Graph Representation:* In a directed graph,  $G = \{V, E\}$  where  $V = \{v_1, v_2, \dots, v_N\}$  represents the set of nodes and  $|E|$  signifies the total number of edges connecting these nodes. The structure of the graph can be illustrated using an adjacency matrix  $A \in \mathbb{R}^{N \times N}$ . Every node within the graph is associated with  $F_N$  features, and the matrix encapsulating these node features is expressed as  $X \in \mathbb{R}^{N \times F_N}$ . A combinatorial Laplacian matrix, denoted as  $L$ , is derived through Equation (1). This involves the use of the degree matrix of  $A$ , symbolised as  $D$ , which is calculated using  $D_{ii} = \sum_{j=1}^N A_{ij}$ .

$$L = I_N - D^{-1/2}AD^{-1/2} \quad (1)$$

2) *Spectral Graph Filtering:* The eigenvectors of the graph Laplacian matrix can be expressed in the Fourier mode as  $\{u_l\}_{l=0}^{N-1} \in \mathbb{R}^N$ , with the Fourier basis  $U = [u_0, \dots, u_{N-1}] \in \mathbb{R}^{N \times N}$ . The corresponding eigenvalues, denoted as  $\{\lambda_l\}_{l=0}^{N-1} \in \mathbb{R}$ , represent the graph Fourier frequencies, and the diagonal matrix containing these Fourier frequencies,  $\Lambda$ , is defined as  $\Lambda = \text{diag}[\lambda_0, \dots, \lambda_{N-1}] \in \mathbb{R}^{N \times N}$ . A signal  $x$  can undergo a graph Fourier transform to become  $\hat{x} = U^T x$ , and the inverse Fourier transform is obtained with  $x = U\hat{x}$ . The convolution operation on the graph  $G$  is defined as:

$$x *_G g = U((U^T x) \odot (U^T g)) \quad (2)$$

where  $g \in \mathbb{R}^N$  denotes a convolutional filter. With  $g_\theta(\Lambda) = \text{diag}(\theta)$ , where  $\theta \in \mathbb{R}^N$  symbolises the vector of Fourier coefficients, the graph convolution of the signal  $x$  is executed as:

$$x *_G g_\theta = U g_\theta(\Lambda) U^T x \quad (3)$$

Given the non-parametric and non-localised nature of the  $g_\theta$  filter, its computational demand is excessively high. Utilising the Chebyshev graph convolution technique, the computational complexity is reduced from  $O(N^2)$  to  $O(KN)$ . The  $g_\theta$  approximation, up to the  $K^{\text{th}}$  order within the Chebyshev polynomial framework, is facilitated using Equation 4. The normalisation of the  $\Lambda$  can be achieved using Equation 5. The term  $\theta_k$  denotes the coefficients of the Chebyshev polynomial, and  $T_k(\hat{\Lambda})$  is derived using Equation 6.

$$g_\theta(\Lambda) = \sum_{k=0}^{K-1} \theta_k T_k(\hat{\Lambda}) \quad (4)$$

$$\hat{\Lambda} = \frac{2\Lambda}{\Lambda_{max}} - I_N \quad (5)$$

$$\{T_0(\hat{\Lambda}) = 1, T_1 = \hat{\Lambda}, T_k(\hat{\Lambda}) = 2\hat{\Lambda}T_{k-1}(\hat{\Lambda}) - T_{k-2}(\hat{\Lambda})\} \quad (6)$$

Ultimately, the graph convolution operation on the signal  $x$  is executed as shown in Equation 7, utilising the normalised Laplacian matrix,  $\tilde{L}$  which is calculated through Equation 8.

$$x *_G g_\theta = U \sum_{k=0}^{K-1} \theta_k T_k(\hat{\Lambda}) U^T x = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{L}) x \quad (7)$$

$$\tilde{L} = \frac{2L}{\lambda_{max}} - I_N \quad (8)$$

3) *Training EEG\_GLT-Net:* In the EEG\_GLT-Net study [18], the classification of EEG MI signals,  $X$  is facilitated through a forward pass using the Spectral GNN function, denoted as  $f(\cdot, \Theta)$ , with a given graph  $G = \{A, X\}$ . The adjacency matrix,  $A$ , integrates  $A_{original}$  and  $m_g$  as outlined in Equation 9. The matrix  $A_{original}$ , defined as  $A_{original_{ij}} = \{0, \text{if } i = j; 1, \text{otherwise}\}$ , is fixed and not subject to training, structured in the dimension of  $\mathbb{R}^{64 \times 64}$ . Meanwhile, the adjacency matrix mask  $m_g \in \mathbb{R}^{64 \times 64}$  is designated as trainable.

$$A = A_{original} \odot m_g \quad (9)$$

EEG MI signals from individual subjects, recorded between  $t = 1s$  and  $t = 3s$ , are trained using Algorithm 1. The detailed structure of the EEG\_GLT-Net is depicted in Figure 1 and Table I, with the specific hyperparameter configurations for the training outlined in Table II. The optimally trained GNN weights ( $\Theta$ ) and the trained adjacency matrix mask ( $m_g$ ) are recorded across various adjacency matrix density levels, ranging from 100% to 13.39%.

4) *EEG MI Time Points GNN Features:* From the pre-trained GNN weights and optimal adjacency matrices across varying  $m_g$  densities ranging from 100% to 13.39%, the set corresponding to a density of 13.39% was chosen. This density was used to extract graph features from EEG MI signals at specific time points, due to its computation efficiency and superior accuracy compared to the 100% set. GNN features were then extracted for all EEG MI time points, spanning

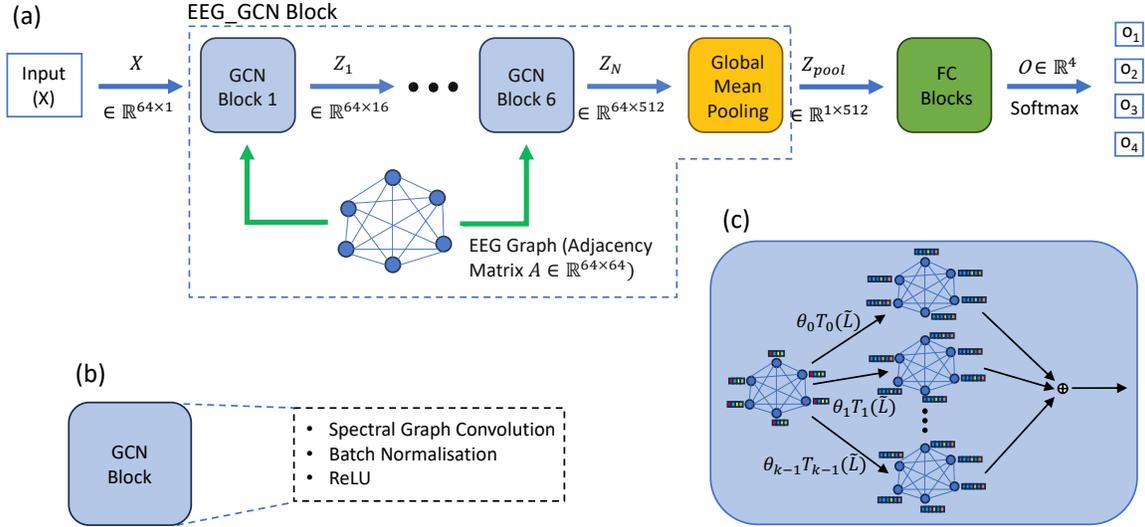


Fig. 1: EEG\_GLT-Net model [18]: (a) Overall architecture (classifying EEG MI of one time point  $\frac{1}{160}$  s of signals from 64 EEG electrodes), (b) Components inside the spectral graph convolution block, (c) Chebyshev spectral graph convolution

TABLE I: Details of EEG\_GLT-Net Model

Layer	Type	Input Size	Polynomial Order	Weights	Output
Input	Input	$64 \times 1$	-	-	-
GCN Blocks					
GC1	Graph Convolution	$64 \times 1$	5	$1 \times 16 \times 5$	$64 \times 16$
BNC1	Batch Normalisation	$64 \times 16$	-	16	$64 \times 16$
GC2	Graph Convolution	$64 \times 16$	5	$16 \times 32 \times 5$	$64 \times 32$
BNC2	Batch Normalisation	$64 \times 32$	-	32	$64 \times 32$
GC3	Graph Convolution	$64 \times 32$	5	$32 \times 64 \times 5$	$64 \times 64$
BNC3	Batch Normalisation	$64 \times 64$	-	64	$64 \times 64$
GC4	Graph Convolution	$64 \times 64$	5	$64 \times 128 \times 5$	$64 \times 128$
BNC4	Batch Normalisation	$64 \times 128$	-	128	$64 \times 128$
GC5	Graph Convolution	$64 \times 128$	5	$128 \times 256 \times 5$	$64 \times 256$
BNC5	Batch Normalisation	$64 \times 256$	-	256	$64 \times 256$
GC6	Graph Convolution	$64 \times 256$	5	$256 \times 512 \times 5$	$64 \times 512$
BNC6	Batch Normalisation	$64 \times 512$	-	512	$64 \times 512$
Global Mean Pooling Block					
P	Global Mean Pool	$64 \times 512$	-	-	512
Fully Connected Blocks					
FC1	Fully Connected	512	-	$512 \times 1024$	1024
BNFC1	Batch Normalisation	1024	-	1024	1024
FC2	Fully Connected	1024	-	$1024 \times 2048$	2048
BNFC2	Batch Normalisation	2048	-	2048	2048
FC3	Fully Connected	$2048 \times 4$	-	$2048 \times 4$	4
S	Softmax Classification	4	-	-	4

TABLE II: Hyperparameter Configuration for Training the EEG\_GLT-Net

Hyperparameter	Value
Training Epochs ( $N_{ep}$ )	1000
Batch Size ( $B$ )	1024
Dropout Rate	0.5
Optimiser	Adam
Initial Learning Rate ( $\eta$ )	0.01

from  $t = 0s$  to  $t = 4s$  for all 84 trials of each subject, was conducted. The GNN feature corresponding to each time point had a dimensionality of  $\mathbb{R}^{512}$ .

#### D. Problem Redefinition

The EEG\_GLT-Net underwent training for the classification of EEG MI time-point signals. Integration the GNN and an optimally trained adjacency matrix significantly enhanced the classification accuracy compared to traditional PCC adjacency matrix method. Nonetheless, ambiguities in signal clarity between different classes at certain time points could adversely affect the model accuracy. Leveraging the high efficacy of the EEG\_GLT-Net model, the pre-trained weights from the GNN and adjacency matrix components were integrated with an RL (Reinforcement Learning) block, resulting in the formation of the EEG\_RL-Net, as depicted in Figure 2.

A reinforcement learning approach is used to train an RL agent for classifying EEG MI time-point signals. Beyond the

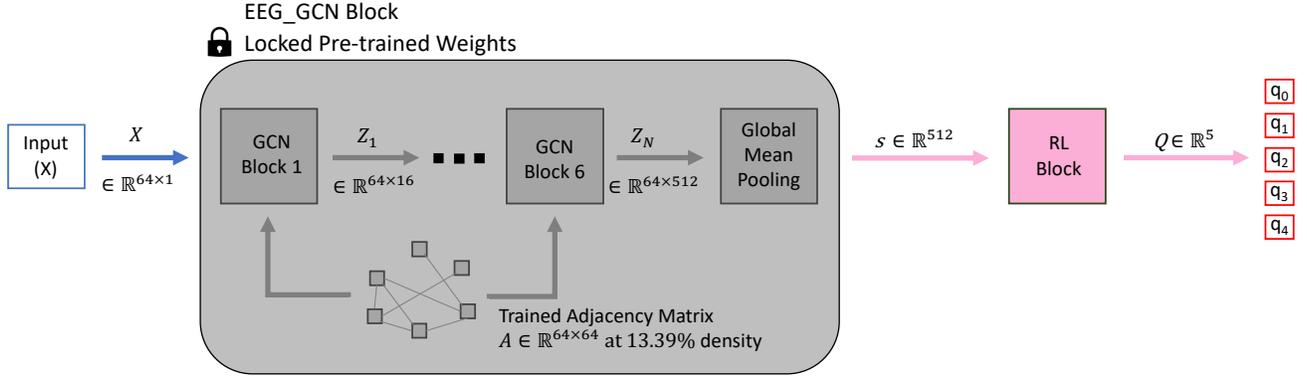


Fig. 2: Overview of the EEG\_RL-Net model: Incorporation of the pre-trained EEG\_GCN Block at a 13.39%  $m_g$  density from the EEG\_GLT-Net, coupled with an RL Block

**Algorithm 1** Finding Optimal EEG\_GCN Weights ( $\Theta$ ) and Adjacency Matrix ( $m_g$ ) at Different Density Levels

**Input:** Graph  $G = \{A, X\}$ , GNN  $f(G, \Theta)$ , GNN initialisation  $\Theta_0$ ,  $A_{original,ij} = \{0, \text{if } i = j; 1, \text{otherwise}\}$ , initial Adjacency Matrix Mask  $m_g^0 = A_{original}$ , learning rate  $\eta = 0.01$ , pruning rate  $p_g = 10\%$ , pre-defined lowest Graph Density Level  $s_g = 13.39\%$ .  
**Output:** Optimal EEG\_GCN weights ( $\Theta^s$ ) with optimal adjacency matrix mask ( $m_g^s$ ) at different graph density levels.

- 1: **while**  $\frac{\|m_g^s\|_0}{\|A_{original}\|_0} \geq s_g$  **do**
- 2:   **for** for iteration  $i = 0, 1, 2, \dots, N_{ep}$  **do**
- 3:     Forward  $f(\cdot, \Theta_i^s)$  with  $G_s = \{m_g^{s,i} \odot A_{original}, X\}$  to compute Cross-Entropy Loss,  $L$
- 4:     Backpropagate and update,  $\Theta_i^s$  and  $m_g^{s,i}$  using Adam Optimiser
- 5:   **end for**
- 6:   Record  $m_g^{s,i}$  and  $\Theta_i^s$  with the highest accuracy in validation set during the  $N_{ep}$  iteration
- 7:   Set  $p_g = 10\%$  of the lowest absolute magnitude values in  $m_g^s$  to 0 and the others to 1, then obtain a new  $m_g^{s+1,0}$
- 8: **end while**

**Algorithm 2** EEG\_RL Environment

- 1: **function** STEP( $s_t, a_t, y_t, s'_t$ )
- 2:   **if**  $a_t = 0$  **then**
- 3:      $r_t = -0.1$
- 4:     Return( $s'_t, r_t$ )
- 5:   **else**
- 6:      $r_t = \begin{cases} r_{right}, \text{ eg. } +10 & \text{if } a_t = y_t \\ r_{wrong}, \text{ eg. } -10 & \text{if } a_t \neq y_t \end{cases}$
- 7:     Return ( $s'_t = Terminal, r_t$ )
- 8:   **end if**
- 9: **end function**

four initial classes, the RL agent has the capability to defer classification of a current time point if it determines that it is not ready. In each state  $s_t$ , the RL agent can perform one of five discrete actions  $a_t \in \{0, 1, 2, 3, 4\}$  within the EEG\_RL

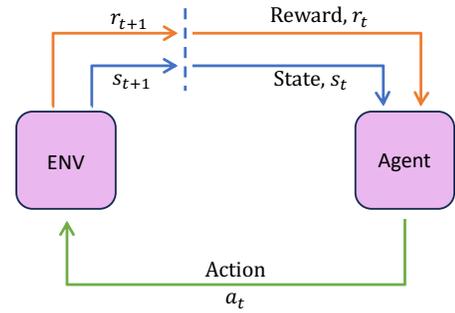


Fig. 3: Agent interaction with EEG\_RL Environment

environment, guided by the GNN features extracted from  $s_t$ . The actions  $a_t$  are described as follows:

- $a_t = 0$  : Skip the current state  $s_t$
- $a_t = 1$  : Classify the signal as Class 1
- $a_t = 2$  : Classify the signal as Class 2
- $a_t = 3$  : Classify the signal as Class 3
- $a_t = 4$  : Classify the signal as Class 4

Following action  $a_t$ , the RL agent is rewarded with  $r_t$  and transitions to the next state  $s'_t$ , as illustrated in Figure 3. Choosing  $a_t = 0$  indicates the agent's hesitance to classify due to uncertainty, leading to a decision to skip the current state with a minimal penalty until it is deemed ready to classify or the episode ends. Upon selecting an action  $a_t > 0$ ,  $s'_t$  is marked as *Terminal*, which concludes the episode and the agent receives  $r_t$ , a positive reward ( $r_{right}$ ) for correct classification or a negative reward ( $r_{wrong}$ ) for incorrect classification. The dynamics of the EEG\_RL environment are elaborated in Algorithm 2. The ultimate goal is for the RL agent to accurately classify EEG MI signals within the designated horizon  $H = 20$  (120 milliseconds) as swiftly as possible.

*E. Data Preprocessing and Data Splitting*

The EEG\_RL-Net training also utilised the PhysioNet dataset, consistent with the approach for EEG\_GLT-Net training. For this training, the entire duration of the EEG MI signals was included, spanning four seconds at a sampling

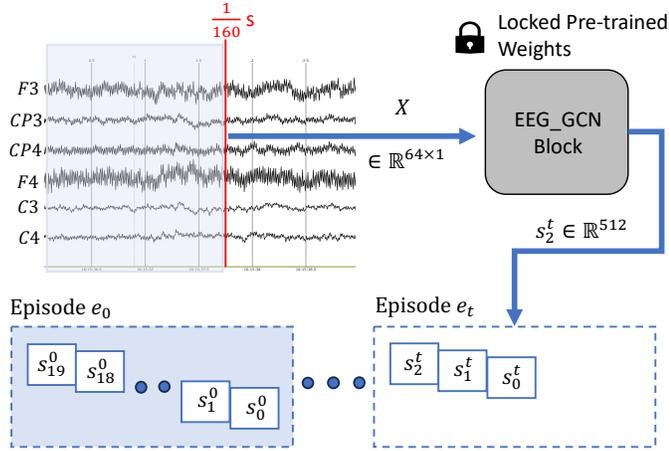


Fig. 4: Conversion of EEG MI time points into states using the pre-trained EEG\_GCN Block, grouped into episodes comprising 20 states each

rate of 160Hz, was included. As outlined in Section II-C4, GNN features of EEG MI time points were used, which were extracted by leveraging pre-trained weights at a 13.39%. The GNN features of each time point were considered as states,  $s \in \mathbb{R}^{512}$ . For all 82 trials, from  $t = 0s$  to  $t = 4s$ , groups of consecutive  $H = 20$  states were organised into episodes without time point overlap between subsequent episodes, creating an episode set,  $E \in \{e_0, e_1, \dots, e_n\}$ , as illustrated in Figure 4. Episodes were randomised using a specific seed, with 80% of the episodes in  $E$  allocated for the training set ( $E_{train}$ ), 10% for the validation set ( $E_{val}$ ), and the remaining 10% for the test set ( $E_{test}$ ).

### F. Dueling Deep Q-Learning

The Deep Q Learning Network (DQN) method, a value-based RL approach, was employed in this study to learn an optimal policy to enable more accurate classification of EEG MI signals. A state-action value,  $Q(s, a)$ , represents the expected discounted reward when the agent is in state  $s$ , and takes action  $a$  according to policy  $\pi$ . With the optimal policy ( $\pi^*$ ), the agent aims to achieve the maximum expected discounted reward  $Q^*(s, a)$ , fulfilling the Bellman equation:

$$Q^*(s, a) = \mathbb{E}_{\pi^*} [r + \gamma \max_{a'} Q^*(s', a') | s, a] \quad (10)$$

here  $r$  is the immediate reward, and  $\gamma$  is the discount factor. The state-action value,  $\hat{Q}(s, a)$ , for state  $s$  and action  $a$  can be approximated using a deep neural network parameterised by  $\theta$ . The loss function is defined as:

$$Loss(\theta) = (\hat{y}^{DQN} - \hat{Q}(s, a; \theta))^2 \quad (11)$$

where  $\hat{y}^{DQN}$  is the target value, calculated as follows:

$$\hat{y}^{DQN} = \begin{cases} r_t, & \text{if } s'_t \text{ is Terminal} \\ r_t + \gamma \max_{a'_t} \hat{Q}(s'_t, a'_t; \theta_{target}), & \text{otherwise} \end{cases} \quad (12)$$

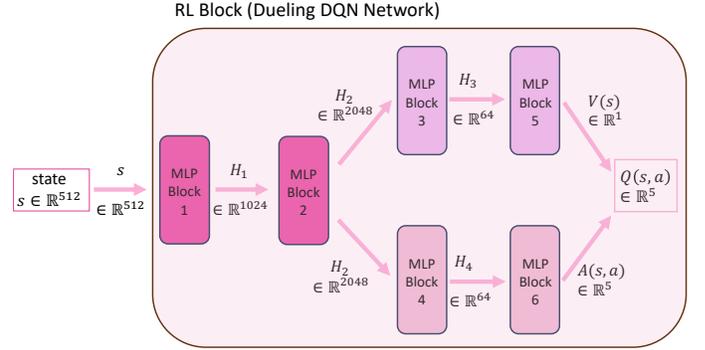


Fig. 5: EEG\_RL-Net's RL Block: Featuring the Dueling Deep Q Network (DQN), this component predicts the q-values linked to various actions

The  $\theta_{target}$  denotes the parameters of the target network, which are kept constant. The approximation  $\hat{Q}(s, a; \theta)$  shares the architecture with the target network. Our study utilises Dueling DQN, a variant of DQN that enhances training stability and efficiency by separating the estimation of  $\hat{Q}(s, a; \theta)$  into state values  $V(s)$  and action advantages  $A(s, a)$ , as follows:

$$\hat{Q}(s, a; \theta) = \hat{V}(s; \alpha) + \hat{A}(s, a; \beta) \quad (13)$$

The network separately estimates the state values and action advantages, which then converge into a single output. The parameters  $\theta$  represent the overall network parameters, with  $\alpha$  and  $\beta$  specifically used for estimating state values and action advantages, respectively. To enhance stability, the equation subtracts the average advantage values from  $\hat{Q}(s, a; \theta)$ :

$$\hat{Q}(s, a; \theta) = \hat{V}(s; \alpha) + [\hat{A}(s, a; \beta) - \frac{1}{|A|} \sum_a \hat{A}(s, a; \beta)] \quad (14)$$

### G. EEG\_RL Algorithm

To generate training data for the RL Block, all possible actions  $a_t = \{0, 1, 2, 3, 4\}$  are executed at each state  $s_t$  within an episode  $e_i$  in  $E_{train}$ , interacting with the EEG\_RL environment to determine the reward  $r_t$  and the subsequent state  $s'_t$ . Each transition records a tuple  $(s, a, r, s')$ . This study employs the Dueling DQN method for the RL block, as illustrated in Figure 5. The Dueling DQN agent undergoes training according to the procedure outlined in Algorithm 3, utilising the Adam optimiser until convergence is achieved. The configuration of the entire EEG\_RL-Net model is outlined in Table III. The parameters of the fixed target network,  $\theta_{target}$ , for the Dueling DQN network, are refreshed after every 50 batch updates of  $\theta$ .

Performance evaluation of the RL agent on  $E_{val}$  and  $E_{test}$  is described in Algorithm 4. At every time step, the agent selects an action based on the q-values predicted by the EEG\_RL-Net. In this study, a correct classification by the agent yields a reward  $r_{right}$ , while an incorrect classification results in  $r_{wrong}$ . The agent's objective in each episode  $e_i$  is to maximise the cumulative reward  $r_{sum}$  within the predefined horizon

**Algorithm 3** Training EEG\_RL-Net’s Dueling DQN Agent

- 1: Initialise randomly Dueling DQN network parameter ( $\theta$ ) and target network parameter ( $\theta_{target}$ ).
- 2: Set of train episodes  $E_{TRAIN} \in \{e_0, e_1, \dots, e_N\}$  where each  $e_i$  has set of states,  $S = \{s_0, s_1, \dots, s_{H-1}\}$ . Each state,  $s_t \in \mathbb{R}^{512}$ .
- 3: At each state  $s_t$ , simulate one step with all possible actions from action set,  $A \in \{0, 1, 2, 3, 4\}$  to observe next state,  $s'_t$  and reward,  $r_t$ . Record all the  $(s_t, a_t, r_t, s'_t)$  tuples to the Buffer  $B$ .
- 4: Shuffle the state transitions in the  $B$  using random seed, and group into mini-batches in size of 64 transitions.
- 5: **for**  $epoch = 0$  to  $EPOCHS$  **do**
- 6:   Compute  $\hat{y}^{DQN}$  for each mini-batch:
- 7:    $\hat{y}^{DQN} = \begin{cases} r_t, & \text{if } s'_t \text{ is Terminal} \\ r_t + \gamma \max_{a'_t} \hat{Q}(s'_t, a'_t; \theta_{target}) & \text{otherwise} \end{cases}$
- 8:    $Loss(\theta) = (\hat{y}^{DQN} - \hat{Q}(s_t, a_t; \theta))^2$
- 9:   Backpropagate to update  $\theta$  using *Adam* optimiser
- 10:   Update  $\theta_{target} = \theta$  at every 50 updates of  $\theta$
- 11: **end for**

$H = 20$ . This requires the agent to make classifications as quickly as possible, since it incurs a penalty of  $r = -0.1$  for each skipped step. However, at time  $t = H - 1$ , skipping is no longer an option, and the agent must make a classification action.

**Algorithm 4** Evaluation of DQN Agent for a Validation or Test Episode

- 1: Episode,  $e_i$  has horizon of  $H = 20$
- 2: At  $e_i$ , the set of states  $S = \{s_0, s_1, \dots, s_{H-1}\}$ , where each  $s_t \in \mathbb{R}^{512}$
- 3: At  $e_i$ , the set of labels  $Y = \{y_0, y_1, \dots, y_{H-1}\}$ , where each  $y_t \in \{1, 2, 3, 4\}$
- 4: Action  $a' \in \{0, 1, 2, 3, 4\}$ , and  $a'' \in \{1, 2, 3, 4\}$
- 5: Initialise  $t = 0$ ,  $r_{sum} = 0$
- 6: **while**  $t < H$  **do**
- 7:    $a_t = \begin{cases} \text{argmax}_{a'} \hat{q}_{DQN}(s_t, a'), & \text{if } t < H - 1 \\ \text{argmax}_{a''} \hat{q}_{DQN}(s_t, a''), & \text{otherwise} \end{cases}$
- 8:    $s'_t, r_t = STEP(s_t, a_t, y_t, s'_t)$
- 9:    $r_{sum} \leftarrow r_{sum} + r_t$
- 10:   **if**  $r_t = Terminate$  **then**
- 11:     Terminate the Episode,  $e_i$
- 12:   **else**
- 13:      $t \leftarrow t + 1$
- 14:   **end if**
- 15: **end while**

*H. Model Setting and Evaluation Metrics*

The structure of EEG\_RL-Net is defined by two principal components: the spectral EEG\_GCN block, which extracts graph features from EEG MI time point signals using pre-trained weights, and the RL block, embodied by the Dueling DQN network. The specifics of the EEG\_RL-Net’s design are provided in Table III. The RL block comprises six MLP

(Multi-Layer Perceptron) layers, or Fully Connected Layers, each followed by a Rectified Linear Unit (ReLU) layer, as described in Equation 15. Information on the training hyperparameters is presented in Table IV. The performance of the different methods was evaluated using both accuracy and F1 score metrics.

TABLE III: Details of EEG\_RL-Net Model

Layer	Type	Input Size	Weights	Output
Input	Input	$64 \times 1$	-	-
EEG_GCN Block				
EEG_GCN	Graph Convolution and Global Pooling	$64 \times 1$	-	512
	RL Block (Dueling DQN Network)			
MLP1	Fully Connected	512	$512 \times 1024$	1024
MLP2	Fully Connected	1024	$1024 \times 2048$	2048
MLP3	Fully Connected	2048	$2048 \times 64$	64
MLP4	Fully Connected	64	$64 \times 1$	1
MLP5	Fully Connected	2048	$2048 \times 64$	64
MLP6	Fully Connected	64	$64 \times 5$	5
Q	Dueling DQN	$64 \times 1$ & $64 \times 5$	-	5

TABLE IV: Hyperparameter Configuration for Training the EEG\_RL-Net

Hyperparameter	Value
Reward Right ( $r_{right}$ )	+10
Reward Wrong ( $r_{wrong}$ )	-10
Reward Skip ( $r_{skip}$ )	-0.1
Discount Factor ( $\gamma$ )	0.99
Training Epoch ( $EPOCHS$ )	150
Batch Size	63
Target Network Update Frequency	50
Initial Learning Rate ( $\eta$ )	0.0001
L2 Regularisation Rate ( $\lambda$ )	0.001
Optimiser	Adam

$$ReLU(x) = \max(0, x) \tag{15}$$

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \tag{16}$$

$$Sensitivity = \frac{TP}{TP + FN} \tag{17}$$

$$Precision = \frac{TP}{TP + FP} \tag{18}$$

$$F1\ Score = \frac{2 \times Precision \times Sensitivity}{Precision + Sensitivity} \tag{19}$$

III. RESULTS AND DISCUSSION

A. EEG\_RL-Net vs EEG\_GLT-Net

Table V shows comparative analysis of mean accuracy between the EEG\_GLT-Net and the EEG\_RL-Net. The EEG\_GLT-Net incorporates two adjacency matrix types: the Pearson Coefficient Correlation (PCC) and the  $m_{g\_GLT}$ . The latter is identified as the most optimal adjacency matrix, after searching through 100% to 13.39% of adjacency matrix density using the EEG\_GLT algorithm. According to paper [18],

employing the  $m_{g\_GLT}$  adjacency matrix yields an accuracy improvement ranging between 0.51% and 22.04% over the PCC adjacency matrix, with significant enhancements noted for subjects  $S_1$  and  $S_{12}$ , at 22.04% and 21.62% respectively. Despite seeing notable improvements in accuracy and F1 score with the  $m_{g\_GLT}$  matrix, certain subjects, specifically  $S_5$ ,  $S_6$ ,  $S_7$ ,  $S_{13}$ ,  $S_{15}$ , and  $S_{19}$ , exhibited classification accuracies below 70%.

TABLE V: Accuracy Assessment: EEG\_RL-Net versus EEG\_GLT-Net

Subj	Accuracy (Mean±Std)		
	EEG_GLT-Net (PCC Adj)	EEG_GLT-Net ( $m_{g\_GLT}$ Adj)	EEG_RL-Net* (our method)
$S_1$	76.47% ± 9.94%	98.51% ± 0.77%	<b>100.00% ± 0.00%</b>
$S_2$	69.13% ± 7.05%	76.18% ± 5.53%	<b>97.73% ± 0.20%</b>
$S_3$	87.28% ± 9.19%	99.17% ± 0.32%	<b>100.00% ± 0.00%</b>
$S_4$	99.13% ± 1.01%	99.97% ± 0.06%	<b>100.00% ± 0.00%</b>
$S_5$	43.19% ± 3.03%	50.95% ± 3.80%	<b>87.72% ± 0.70%</b>
$S_6$	58.23% ± 5.19%	69.60% ± 5.67%	<b>90.89% ± 1.50%</b>
$S_7$	50.98% ± 3.80%	59.45% ± 3.00%	<b>89.24% ± 2.10%</b>
$S_8$	95.06% ± 5.96%	99.95% ± 0.07%	<b>100.00% ± 0.00%</b>
$S_9$	97.64% ± 3.33%	99.95% ± 0.08%	<b>100.00% ± 0.00%</b>
$S_{10}$	99.24% ± 0.19%	99.99% ± 0.01%	<b>100.00% ± 0.00%</b>
$S_{11}$	99.48% ± 0.70%	99.99% ± 0.01%	<b>100.00% ± 0.00%</b>
$S_{12}$	78.07% ± 8.95%	99.69% ± 0.32%	<b>100.00% ± 0.00%</b>
$S_{13}$	41.35% ± 1.23%	44.50% ± 2.23%	<b>89.45% ± 0.90%</b>
$S_{14}$	55.97% ± 6.47%	72.39% ± 6.43%	<b>91.59% ± 2.10%</b>
$S_{15}$	52.11% ± 3.96%	67.55% ± 9.26%	<b>80.83% ± 1.50%</b>
$S_{16}$	96.75% ± 5.00%	99.98% ± 0.03%	<b>100.00% ± 0.00%</b>
$S_{17}$	98.83% ± 2.33%	99.98% ± 0.03%	<b>100.00% ± 0.00%</b>
$S_{18}$	86.19% ± 9.95%	99.92% ± 0.12%	<b>100.00% ± 0.00%</b>
$S_{19}$	38.38% ± 2.27%	41.41% ± 1.44%	<b>79.65% ± 1.40%</b>
$S_{20}$	98.44% ± 0.68%	99.94% ± 0.11%	<b>100.00% ± 0.00%</b>
Overall	76.10% ± 22.71%	83.95% ± 21.43%	<b>95.36% ± 6.83%</b>

\*  $r_{right} = +10$ ,  $r_{wrong} = -10$ ,  $r_{skip} = -0.1$ ,  $H = 20$

Using baseline parameters ( $r_{right} = +10$ ,  $r_{wrong} = -10$ ,  $r_{skip} = -0.1$  and  $H = 20$ ), the EEG\_RL-Net framework advances the accuracy beyond the current state-of-the-art EEG\_GLT-Net employing the  $m_{g\_GLT}$  adjacency matrix, with improvements spanning 0.01% to 44.95%. A total of 12 out of 20 subjects, namely  $S_1$ ,  $S_3$ ,  $S_4$ ,  $S_8$ ,  $S_9$ ,  $S_{10}$ ,  $S_{11}$ ,  $S_{12}$ ,  $S_{16}$ ,  $S_{17}$ ,  $S_{18}$ , and  $S_{20}$ , achieved perfect classification. The EEG\_RL-Net also significantly elevated the accuracies for  $S_{13}$  and  $S_{19}$  to 89.45% and 79.65%, respectively. Even for subjects  $S_{13}$  and  $S_{19}$ , who initially demonstrated low accuracies, modest improvement in accuracy at 44.50% and 41.41%, respectively was achieved using the EEG\_GLT-Net with the  $m_{g\_GLT}$  matrix.

The EEG\_GLT-Net with the  $m_{g\_GLT}$  matrix boosted accuracy across the 20 subjects, increasing the average accuracy by 7.85% (from 76.10% to 83.95%). Given the inherent noise in EEG MI time-point signals and the challenge of classifying signals representing  $\frac{1}{160}$  s, the EEG\_GLT-Net showed a decline in performance accuracy due to its attempt to classify all time points. Comparatively, the EEG\_RL-Net achieved remarkable increase in average accuracy across the 20 subjects to 95.35%. This substantial improvement is the result of the RL agent’s capacity to discern the appropriateness of the current signal for classification. The agent has been optimised to classify signals as swiftly as possible within a 20 time-point window, averaging a classification time of 2.91 time points in the

EEG\_RL-Net setup.

### B. Study of Changing $r_{right}$ Values

Table VI demonstrates the effect of varying the  $r_{right}$  value (+5, +10, +15), on average accuracy while keeping  $r_{wrong} = -10$  constant. The results show average accuracies of 95.57%, 95.36%, and 94.94% for  $r_{right} = +5$ , +10, and +20, respectively. Notably, the accuracy tends to improve when  $r_{right}$  is less than  $r_{wrong}$ , but declines when  $r_{right}$  exceeds  $r_{wrong}$ , although the level of variance is minimal at just 0.63%.

TABLE VI: Impact of Varying  $r_{right}$  Values on Accuracy and Classification Time

Subj	Mean Accuracy (Mean Classification Time)		
	$r_r = +5$ $r_w = -10$	$r_r = +10$ $r_w = -10$	$r_r = +20$ $r_w = -10$
$S_1$	99.93% (1.70)	100.00% (1.80)	100.00% (1.50)
$S_2$	97.86% (1.51)	97.73% (1.87)	97.86% (1.65)
$S_3$	100.00% (2.10)	100.00% (2.20)	100.00% (1.90)
$S_4$	100.00% (2.80)	100.00% (2.50)	100.00% (1.80)
$S_5$	87.65% (5.71)	87.72% (4.55)	86.14% (3.45)
$S_6$	91.10% (3.37)	90.89% (2.80)	89.72% (2.07)
$S_7$	91.24% (3.64)	89.24% (3.23)	87.80% (3.66)
$S_8$	100.00% (2.90)	100.00% (2.80)	100.00% (2.10)
$S_9$	100.00% (3.70)	100.00% (3.00)	100.00% (2.20)
$S_{10}$	99.93% (2.40)	100.00% (2.20)	100.00% (1.80)
$S_{11}$	100.00% (2.00)	100.00% (1.50)	100.00% (2.30)
$S_{12}$	100.00% (2.80)	100.00% (2.40)	100.00% (2.60)
$S_{13}$	89.59% (4.90)	89.45% (3.97)	88.28% (3.58)
$S_{14}$	93.45% (3.75)	91.59% (3.14)	89.86% (2.81)
$S_{15}$	82.89% (5.23)	80.83% (4.80)	79.45% (4.51)
$S_{16}$	100.00% (2.40)	100.00% (2.70)	100.00% (2.20)
$S_{17}$	100.00% (3.00)	100.00% (2.00)	100.00% (1.90)
$S_{18}$	100.00% (3.30)	100.00% (2.40)	100.00% (1.50)
$S_{19}$	77.79% (6.89)	79.65% (5.64)	79.59% (4.80)
$S_{20}$	100.00% (2.20)	100.00% (2.70)	100.00% (2.10)
Mean	95.57% (3.32)	95.36% (2.91)	94.94% (2.51)
Std	± 6.72%	± 6.83%	± 7.32%

$r_{skip} = -0.1$  and  $H = 20$

On an individual basis,  $r_{right} = +5$  yielded higher accuracies for most subjects, except for  $S_1$ ,  $S_5$ , and  $S_{19}$ , where  $r_{right} = +10$  performed marginally better. No subjects showed improved accuracy when  $r_{right}$  was greater than  $r_{wrong}$ . Therefore, for optimal performance, the magnitude of  $r_{right}$  should not exceed  $r_{wrong}$ . It appears that accuracy is enhanced by a higher penalty for incorrect classifications ( $r_{wrong}$ ) rather than a higher reward for correct ones ( $r_{right}$ ) enhances accuracy, likely motivating the agent to avoid misclassifications more stringently.

Regarding the time points required to classify EEG MI signals, the configuration with  $r_{right} = +10$  and  $r_{wrong} = -10$  averages at 2.91 time points. Increasing  $r_{right}$  to +20 (while  $r_{wrong}$  remains at -10) reduces the classification time to 2.51 time points. Conversely, lowering  $r_{right}$  to +5 increases the average classification time to 3.32 time points, indicating a more cautious approach by the agent, likely due to prioritising accuracy over speed by utilising the option to skip uncertain classifications.

### C. Study of Changing $r_{wrong}$ Values

In this study, we examined the impact of altering the  $r_{wrong}$  values while keeping the  $r_{right}$  constant at +10, as shown in

Table VII. We observed the  $r_{wrong}$  values at  $-10, -20, -30,$  and  $-40,$  correlating with an average performance accuracy of  $95.35\%, 95.18\%, 95.11\%,$  and  $94.88\%,$  respectively. This indicates that simply increasing the negative magnitude of  $r_{wrong}$  beyond that of  $r_{right}$  does not invariably lead to enhanced performance accuracy. Additionally, we found that the time required for signal classification was directly related to the difference in rewards.

TABLE VII: Impact of Varying  $r_{wrong}$  Values on Accuracy and Classification Time

Subj	Mean Accuracy (Mean Classification Time)			
	$r_r = +10$	$r_r = +10$	$r_r = +10$	$r_r = +10$
	$r_w = -10$	$r_w = -20$	$r_w = -30$	$r_w = -40$
$S_1$	100.00% (1.80)	99.79% (1.60)	100.00% (2.00)	99.93% (2.00)
$S_2$	97.73% (1.87)	98.21% (1.97)	97.93% (2.49)	97.93% (2.88)
$S_3$	100.00% (2.20)	100.00% (1.70)	100.00% (2.70)	100.00% (2.20)
$S_4$	100.00% (2.50)	100.00% (3.30)	100.00% (4.30)	100.00% (5.70)
$S_5$	87.72% (4.55)	86.90% (6.26)	85.38% (6.28)	85.79% (7.31)
$S_6$	90.89% (2.80)	90.00% (2.63)	90.90% (4.94)	90.41% (4.25)
$S_7$	89.24% (3.23)	89.59% (5.09)	89.31% (5.20)	88.90% (6.70)
$S_8$	100.00% (2.80)	100.00% (2.60)	100.00% (3.70)	100.00% (4.30)
$S_9$	100.00% (3.00)	100.00% (3.30)	100.00% (5.00)	100.00% (5.30)
$S_{10}$	100.00% (2.20)	99.93% (1.90)	99.86% (2.00)	99.93% (2.52)
$S_{11}$	100.00% (1.50)	100.00% (2.20)	100.00% (3.00)	100.00% (3.30)
$S_{12}$	100.00% (2.40)	100.00% (2.70)	100.00% (4.10)	100.00% (5.10)
$S_{13}$	89.45% (3.97)	89.52% (5.33)	89.52% (6.24)	87.38% (6.07)
$S_{14}$	91.59% (3.14)	91.31% (3.00)	92.14% (3.00)	90.55% (4.37)
$S_{15}$	80.83% (4.80)	81.38% (4.53)	80.14% (3.48)	80.69% (4.13)
$S_{16}$	100.00% (2.70)	100.00% (3.10)	100.00% (5.00)	100.00% (4.30)
$S_{17}$	100.00% (2.00)	100.00% (3.70)	100.00% (3.20)	100.00% (4.60)
$S_{18}$	100.00% (2.40)	100.00% (3.20)	100.00% (2.40)	100.00% (4.30)
$S_{19}$	79.65% (5.64)	76.90% (7.21)	76.97% (8.22)	76.14% (9.20)
$S_{20}$	100.00% (2.70)	100.00% (2.80)	100.00% (3.40)	100.00% (4.80)
Mean	95.36% (2.91)	95.18% (3.41)	95.11% (4.03)	94.88% (4.66)
Std	$\pm 6.83\%$	$\pm 7.19\%$	$\pm 7.37\%$	$\pm 7.57\%$

$r_{skip} = -0.1$  and  $H = 20$

Despite the reward configuration of  $\{r_{right} = +10, r_{wrong} = -10\}$  achieving the highest average performance accuracy among the four settings, it does not universally outperform across all test subjects. Specifically, this configuration was only superior for subjects  $S_5$  and  $S_6$ . Conversely, the configuration of  $\{r_{right} = +10, r_{wrong} = -20\}$  exhibited higher performance accuracy in subjects  $S_2, S_7, S_{13},$  and  $S_{15}$ . For subject  $S_{14}$ , the  $\{r_{right} = +10, r_{wrong} = -30\}$  setting was more advantageous.

Although a smaller magnitude of  $r_{wrong}$  relative to  $r_{right}$  appears beneficial, a higher  $r_{wrong}$  to  $r_{right}$  ratio does not necessarily equate to improved accuracy. As demonstrated in Table VII, performance accuracy diminishes with an increasing ratio, identifying the optimal ratio as twice the magnitude of  $r_{wrong}$  to  $r_{right}$ . Furthermore, comparing different of reward configurations with equivalent magnitude ratios, such as  $\{r_{right} = +5, r_{wrong} = -10\}$  and  $\{r_{right} = +10, r_{wrong} = -20\}$ , reveal subtle differences are noted in average performance accuracy and classification time. The former configuration outperforms in both average accuracy and time efficiency for classification.

According to Table VII, the classification time escalates with the increases in  $r_{wrong}$  magnitude, where average times of 2.91, 3.41, 4.03, and 4.66 seconds were recorded for  $r_{wrong}$  values of  $-10, -20, -30,$  and  $-40,$  respectively. This

trend suggests that as the penalty for incorrect classification outweighs the reward for correct answers, the agents proceed with increased caution, hence extending the classification time.

#### D. Effects of Episode Length Variation and Optimisation on Classification Performance

In this study, we examined the mean accuracy, F1 score, and mean classification time across various episode lengths ( $H$ ), including 10, 20, 30, and 40, as presented in Table VIII. We observed that both accuracy and F1 scores increased with extension of the episode horizon extends. Conversely, classification time per point increased with longer episode lengths. These finding suggests that larger episode lengths contribute to improvements in accuracy and F1 scores.

TABLE VIII: Impact of Varying Episode Lengths ( $H$ ) Values on Accuracy, F1 Score and Classification Time

Horizon ( $H$ )	Accuracy (Mean $\pm$ Std)	F1 Score (Mean $\pm$ Std)	Mean Classification Time
10	94.46% $\pm$ 8.10%	94.42% $\pm$ 8.15%	2.18
20	95.14% $\pm$ 7.14%	95.10% $\pm$ 7.18%	3.76
30	95.56% $\pm$ 6.54%	95.53% $\pm$ 5.53%	5.53
40	95.82% $\pm$ 6.16%	95.79% $\pm$ 6.54%	6.54

Table IX delineates the optimal configuration of reward for correct ( $r_{right}$ ) and incorrect ( $r_{wrong}$ ) decisions, and episode horizon ( $H$ ) that achieves the highest accuracy and F1 score in the shortest classification time possible. In this optimal setting, the RL agent demonstrates superior performance, achieving an average accuracy of  $96.40\%$  and an average classification time of less than 25 milliseconds across all 20 subjects.

TABLE IX: Subject-wise Classification Accuracy and Time with Optimal Reward Settings and Episode Lengths

Subj	Mean Accuracy	Mean F1 Score	Mean Classification Time	$(r_{right}, r_{wrong})$	Episode Horizon
$S_1$	100.00%	100.00%	1.45	(20, -10)	10
$S_2$	98.65%	98.62%	2.93	(20, -30)	30
$S_3$	100.00%	100.00%	1.13	(20, -10)	10
$S_4$	100.00%	100.00%	1.32	(20, -30)	10
$S_5$	90.21%	90.05%	4.85	(10, -10)	30
$S_6$	92.06%	92.06%	4.25	(5, -10)	40
$S_7$	92.33%	92.29%	9.94	(10, -30)	40
$S_8$	100.00%	100.00%	1.23	(10, -10)	10
$S_9$	100.00%	100.00%	1.27	(20, -10)	10
$S_{10}$	100.00%	100.00%	1.17	(20, -10)	10
$S_{11}$	100.00%	100.00%	1.11	(20, -20)	10
$S_{12}$	100.00%	100.00%	1.19	(20, -10)	10
$S_{13}$	93.29%	93.27%	5.95	(10, -10)	40
$S_{14}$	93.70%	93.69%	4.17	(5, -10)	40
$S_{15}$	85.48%	85.43%	7.51	(10, -20)	40
$S_{16}$	100.00%	100.00%	1.25	(20, -10)	10
$S_{17}$	100.00%	100.00%	1.28	(20, -10)	10
$S_{18}$	100.00%	100.00%	1.27	(10, -10)	10
$S_{19}$	82.33%	82.20%	9.69	(20, -30)	40
$S_{20}$	100.00%	100.00%	1.15	(20, -10)	10
Mean	96.40%	96.38%	3.21	-	-
Std	$\pm 5.47$	$\pm 5.50$	-	-	-

Our analysis, as indicated in Table IX shows that the RL agent achieved accuracy exceeding  $90.00\%$  for each subject, with the exceptions of  $S_{15}$  and  $S_{19}$  whose accuracies were

85.48% and 82.33%, respectively. Subjects such as  $S_1$ ,  $S_3$ ,  $S_4$ ,  $S_8$ ,  $S_9$ ,  $S_{10}$ ,  $S_{11}$ ,  $S_{12}$ ,  $S_{16}$ ,  $S_{17}$ ,  $S_{18}$ , and  $S_{20}$ , where the RL agent achieved perfect classification, had notably clearer EEG MI signals. For these subjects, the agent performed consistently well across most reward and episode horizon configurations. Classifications were achieved within an average of 2 time points, where an optimal episode horizon of 10 and a reward configuration where  $r_{right}$  significantly exceeded  $r_{wrong}$  were conducive to faster classification decisions.

Particularly noteworthy was the performance of EEG\_RL-Net on subject  $S_{13}$ , where the RL agent achieved a classification accuracy of 93.29%. This represented an exceptional improvement by 48.79% over EEG\_GLT-Net with  $m_{g\_GLT}$ , the current state-of-the-art EEG MI time point classification method. The classification for  $S_{13}$  took 6 time points on average, possibly reflecting the only subtle distinctions between EEG MI tasks for this subject.

#### IV. CONCLUSION

Our study introduces EEG\_RL-Net, an innovative approach for the real-time classification of EEG-based motor imagery (MI) signals utilising reinforcement learning (RL) techniques. Building on the foundation of EEG\_GLT-Net's EEG\_GCN block and optimising computational efficiency with an adjacency matrix density of just 13.39%, EEG\_RL-Net not only achieves accurate classification of EEG MI signals but also identifies signals that are unsuitable for classification. Remarkably, it achieved 100.00% classification accuracy for 12 out of 20 subjects within less than 12.5 milliseconds. For challenging subjects ( $S_{13}$  and  $S_{19}$  in this study), where previous state-of-the-art methods such as EEG\_GLT-Net could classify with accuracies of only 44.50% and 41.41% respectively, EEG\_RL-Net achieved unprecedented improvement in performance, reaching classification accuracies of 93.29% and 82.33% in less than 62.5 milliseconds. These results underscore the robustness and efficacy of EEG\_RL-Net in enhancing classification rates, filling a gap for subjects previously deemed difficult by existing classification methods. In future work, we will further explore the integration of the optimal adjacency matrix  $m_{g\_GLT}$  for advanced graph feature extraction in the EEG\_GCN block, aiming to unlock even greater improvements in the classification capabilities of our EEG\_RL-Net system.

#### REFERENCES

- [1] J. R. Wolpaw, N. Birbaumer, D. J. McFarland, G. Pfurtscheller, and T. M. Vaughan, "Brain-computer interfaces for communication and control," *Clinical Neurophysiology*, vol. 113, no. 6, pp. 767–791, 2002. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1388245702000573>
- [2] M. A. Lebedev and M. A. Nicolelis, "Brain-machine interfaces: past, present and future," *TRENDS in Neurosciences*, vol. 29, no. 9, pp. 536–546, 2006.
- [3] D. L. Schomer and F. L. Da Silva, *Niedermeyer's electroencephalography: basic principles, clinical applications, and related fields*. Lippincott Williams & Wilkins, 2012.
- [4] J. Cao, L. Yang, P. G. Sarrigiannis, D. Blackburn, and Y. Zhao, "Dementia classification using a graph neural network on imaging of effective brain connectivity," *Computers in Biology and Medicine*, vol. 168, p. 107701, 2024.
- [5] Z. Wang, C. Hu, W. Liu, X. Zhou, and X. Zhao, "Eeg-based high-performance depression state recognition," *Frontiers in Neuroscience*, vol. 17, p. 1301214, 2024.
- [6] W. Cappelletti, Y. Xie, and P. Frossard, "Learning self-supervised dynamic networks for seizure analysis," in *ICLR 2024 Workshop on Learning from Time Series For Health*, 2024.
- [7] Y. Ding, N. Robinson, C. Tong, Q. Zeng, and C. Guan, "Lggnnet: Learning from local-global-graph representations for brain-computer interface," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [8] Y. Hou, S. Jia, X. Lun, Z. Hao, Y. Shi, Y. Li, R. Zeng, and J. Lv, "Gens-net: a graph convolutional neural network approach for decoding time-resolved eeg motor imagery signals," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [9] M. A. Abbasi, H. F. Abbasi, M. Z. Aziz, W. Haider, Z. Fan, and X. Yu, "A novel precisely designed compact convolutional eeg classifier for motor imagery classification," *Signal, Image and Video Processing*, pp. 1–12, 2024.
- [10] J. Hubbard, A. Kikumoto, and U. Mayr, "Eeg decoding reveals the strength and temporal dynamics of goal-relevant representations," *Scientific reports*, vol. 9, no. 1, p. 9051, 2019.
- [11] D. J. McFarland, L. A. Miner, T. M. Vaughan, and J. R. Wolpaw, "Mu and beta rhythm topographies during motor imagery and actual movements," *Brain topography*, vol. 12, pp. 177–186, 2000.
- [12] A. Biasiucci, R. Leeb, I. Iturrate, S. Perdakis, A. Al-Khodairy, T. Corbet, A. Schneider, T. Schmidlin, H. Zhang, M. Bassolino *et al.*, "Brain-actuated functional electrical stimulation elicits lasting arm motor recovery after stroke," *Nature communications*, vol. 9, no. 1, p. 2421, 2018.
- [13] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1915–1929, 2012.
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [15] Y. LeCun, K. Kavukcuoglu, and C. Farabet, "Convolutional networks and applications in vision," in *Proceedings of 2010 IEEE international symposium on circuits and systems*. IEEE, 2010, pp. 253–256.
- [16] Z. Zhang, P. Cui, and W. Zhu, "Deep learning on graphs: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp. 249–270, 2020.
- [17] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 1, pp. 4–24, 2020.
- [18] H. W. Aung, J. J. Li, Y. An, and S. W. Su, "Eeg\_glt-net: Optimising eeg graphs for real-time motor imagery signals classification," 2024.
- [19] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.
- [20] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein, "Geometric deep learning on graphs and manifolds using mixture model cnns," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5115–5124.
- [21] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning convolutional neural networks for graphs," in *International conference on machine learning*. PMLR, 2016, pp. 2014–2023.
- [22] H. Gao, Z. Wang, and S. Ji, "Large-scale learnable graph convolutional networks," in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 1416–1424.
- [23] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- [24] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," *Advances in neural information processing systems*, vol. 29, 2016.
- [25] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, "Cayleynets: Graph convolutional neural networks with complex rational spectral filters," *IEEE Transactions on Signal Processing*, vol. 67, no. 1, pp. 97–109, 2018.
- [26] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE signal processing magazine*, vol. 30, no. 3, pp. 83–98, 2013.

- [27] G. Bao, K. Yang, L. Tong, J. Shu, R. Zhang, L. Wang, B. Yan, and Y. Zeng, "Linking multi-layer dynamical gcn with style-based recalibration cnn for eeg-based emotion recognition," *Frontiers in Neurobotics*, vol. 16, p. 834952, 2022.
- [28] D. Zeng, K. Huang, C. Xu, H. Shen, and Z. Chen, "Hierarchy graph convolution network and tree classification for epileptic detection on electroencephalography signals," *IEEE transactions on cognitive and developmental systems*, vol. 13, no. 4, pp. 955–968, 2020.
- [29] L. Meng, J. Hu, Y. Deng, and Y. Hu, "Electrical status epilepticus during sleep electroencephalogram waveform identification and analysis based on a graph convolutional neural network," *Biomedical Signal Processing and Control*, vol. 77, p. 103788, 2022.
- [30] H. Li, H. Ji, J. Yu, J. Li, L. Jin, L. Liu, Z. Bai, and C. Ye, "A sequential learning model with gnn for eeg-emg-based stroke rehabilitation bci," *Frontiers in Neuroscience*, vol. 17, p. 1125230, 2023.
- [31] R. Zhang, Z. Wang, F. Yang, and Y. Liu, "Recognizing the level of organizational commitment based on deep learning methods and eeg," in *ITM Web of Conferences*, vol. 47. EDP Sciences, 2022, p. 02044.
- [32] M. Jia, W. Liu, J. Duan, L. Chen, C. Chen, Q. Wang, and Z. Zhou, "Efficient graph convolutional networks for seizure prediction using scalp eeg," *Frontiers in Neuroscience*, vol. 16, p. 967116, 2022.
- [33] N. Wagh and Y. Varatharajah, "Eeg-gcnn: Augmenting electroencephalogram-based neurological disease diagnosis using a domain-guided graph convolutional neural network," in *Machine Learning for Health*. PMLR, 2020, pp. 367–378.
- [34] N. Khaleghi, T. Y. Rezaii, S. Beheshti, and S. Meshgini, "Developing an efficient functional connectivity-based geometric deep network for automatic eeg-based visual decoding," *Biomedical Signal Processing and Control*, vol. 80, p. 104221, 2023.
- [35] W. Ma, C. Wang, X. Sun, X. Lin, and Y. Wang, "A double-branch graph convolutional network based on individual differences weakening for motor imagery eeg classification," *Biomedical Signal Processing and Control*, vol. 84, p. 104684, 2023.
- [36] T. Song, W. Zheng, P. Song, and Z. Cui, "Eeg emotion recognition using dynamical graph convolutional neural networks," *IEEE Transactions on Affective Computing*, vol. 11, no. 3, pp. 532–541, 2018.
- [37] T. Chen, Y. Sui, X. Chen, A. Zhang, and Z. Wang, "A unified lottery ticket hypothesis for graph neural networks," in *International conference on machine learning*. PMLR, 2021, pp. 1695–1706.
- [38] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "Physiobank, physiotookit, and physionet: components of a new research resource for complex physiologic signals," *circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [39] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [40] J. Janisch, T. Pevný, and V. Lisý, "Classification with costly features using deep reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3959–3966.
- [41] D. Dua and C. Graff, "Uci machine learning repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [42] C. Song, C. Chen, Y. Li, and X. Wu, "Deep reinforcement learning apply in electromyography data classification," in *2018 IEEE International Conference on Cyborg and Bionic Systems (CBS)*. IEEE, 2018, pp. 505–510.