

# E<sup>2</sup>-LLM: Bridging Neural Signals and Interpretable Affective Analysis

Anonymous ACL submission

## Abstract

Emotion recognition from electroencephalography (EEG) signals remains challenging due to high inter-subject variability, limited labeled data, and the lack of interpretable reasoning in existing approaches. While recent multimodal large language models (MLLMs) have advanced emotion analysis, they have not been adapted to handle the unique spatiotemporal characteristics of neural signals. We present E<sup>2</sup>-LLM (EEG-to-Emotion Large Language Model), the first MLLM framework for interpretable emotion analysis from EEG. E<sup>2</sup>-LLM integrates a pretrained EEG encoder with Qwen-based LLMs through learnable projection layers, employing a multi-stage training pipeline that encompasses emotion-discriminative pretraining, cross-modal alignment, and instruction tuning with chain-of-thought reasoning. We design a comprehensive evaluation protocol covering basic emotion prediction, multi-task reasoning, and zero-shot scenario understanding. Experiments on the dataset across seven emotion categories demonstrate that E<sup>2</sup>-LLM achieves excellent performance on emotion classification, with larger variants showing enhanced reliability and superior zero-shot generalization to complex reasoning scenarios. Our work establishes a new paradigm combining physiological signals with LLM reasoning capabilities, showing that model scaling improves both recognition accuracy and interpretable emotional understanding in affective computing.

## 1 Introduction

Multimodal emotion recognition aims to infer human emotional states by integrating heterogeneous signals such as speech, text, facial expressions, and physiological cues (Shou et al., 2025b). By fusing complementary information across modalities, multimodal approaches achieve more robust affective inference than unimodal methods. Deep learning architectures have evolved from CNNs

and RNNs to Transformer-based models that align asynchronous multimodal signals using multi-head self-attention (Yao et al., 2024; Chen et al., 2025). Recent work has explored graph neural networks to model cross-modal dependencies through message passing (Shou et al., 2025a). However, these methods remain limited to fixed emotion taxonomies and struggle with the subjectivity and ambiguity inherent in affective labeling (Wang et al., 2024), particularly when fine-grained semantic reasoning is required.

Electroencephalography (EEG), a non-invasive neuroimaging technique recording brain electrical activity through scalp electrodes, shows substantial promise in brain-computer interfaces, cognitive assessment, and neurological diagnostics (Babu et al., 2025). However, EEG signal processing faces significant challenges including artifact contamination from eye movements and muscle activity, high inter- and intra-subject variability, and limited large-scale annotated datasets (Huang et al., 2023). Meanwhile, Large Language Models (LLMs) such as GPT and BERT have revolutionized natural language processing through transformer architectures that capture sequential dependencies and enable few-shot and zero-shot learning. The self-attention mechanism in LLMs dynamically weights input features and captures long-range dependencies, making it well-suited for modeling EEG’s spatiotemporal dynamics. Recent work has explored convergence of LLMs and EEG analysis, with foundation models like LaBraM (Jiang et al., 2024c) and EEGPT (Wang et al., 2024) using transformer architectures with masked or autoregressive pretraining to learn generalizable EEG representations. For emotion recognition, Multimodal Large Language Models (MLLMs) have achieved remarkable progress by integrating heterogeneous signals through specialized connector modules, enabling end-to-end multimodal emotion reasoning that surpasses traditional models without extensive labeled

085	data (Yang et al., 2025b).	
086	Despite recent advances, existing EEG-based	
087	emotion recognition approaches have several key	
088	limitations. First, while LLM-based methods have	
089	shown promise in EEG-to-text translation (Mishra	
090	et al., 2025) and foundation models (Jiang et al.,	
091	2024a), they lack emotion-specific reasoning capa-	
092	bilities. Recent emotion-focused MLLMs (Yang	
093	et al., 2025b) have not been adapted to handle	
094	the unique spatiotemporal characteristics of neural	
095	signals. Second, conventional systems operate as	
096	closed-set classifiers that output categorical labels	
097	without interpretable rationales, lacking hierarchi-	
098	cal cross-modal alignment mechanisms to capture	
099	both local cues (micro facial expressions, pitch	
100	shifts) and global affective patterns. Third, early	
101	LLM-based approaches rely solely on textual fea-	
102	tures, failing to capture the physiological dynamics	
103	central to emotional states. Although multimodal	
104	emotion benchmarks like EMER (Lian et al., 2024)	
105	and MERR (Liu et al., 2024a) provide reasoning	
106	annotations, no prior work has systematically inte-	
107	grated EEG signals with chain-of-thought emotion	
108	reasoning through end-to-end instruction tuning	
109	to enable interpretable, open-vocabulary emotion	
110	analysis from neural data.	
111	To address these limitations, we propose E <sup>2</sup> -	
112	LLM (EEG-to-Emotion Large Language Model),	
113	the first multimodal large language model for inter-	
114	pretable emotion analysis from EEG signals. E <sup>2</sup> -	
115	LLM features three key innovations: (1) A hier-	
116	archical architecture integrating a pretrained EEG	
117	encoder with Qwen3-based LLMs via learnable	
118	projection layers that align spatiotemporal EEG	
119	dynamics with language embeddings; (2) A multi-	
120	stage training pipeline encompassing emotion-	
121	discriminative pretraining, cross-modal alignment,	
122	and instruction tuning with chain-of-thought rea-	
123	soning to establish semantic correspondence while	
124	enabling interpretable analysis; (3) A compre-	
125	hensive evaluation protocol covering basic emotion	
126	prediction, multi-task reasoning (pairwise compar-	
127	ison, superlative selection, individual matching),	
128	and zero-shot scenario reasoning. Evaluated on the	
129	SEED-VII dataset (Jiang et al., 2024a) across seven	
130	emotion categories, E <sup>2</sup> -LLM obtains excellent per-	
131	formance, with larger variants demonstrating en-	
132	hanced zero-shot generalization to unseen complex	
133	scenarios. Our work establishes a new paradigm	
134	combining physiological signals with LLM rea-	
135	soning, showing that model scaling enhances both	
136	recognition accuracy and interpretable emotional	
	understanding in affective computing.	137
	<b>2 Related Work</b>	138
	<b>LLM-based EEG Analysis</b> The integration of	139
	large language models (LLMs) with electroen-	140
	cephalography (EEG) has emerged as a promis-	141
	ing research direction. Recent advances span	142
	several key areas: foundation models such as	143
	LaBraM (Jiang et al., 2024c) and EEGPT (Wang	144
	et al., 2024) adopt transformer-based architec-	145
	tures with masked or autoregressive pretrain-	146
	ing to learn generalizable EEG representations	147
	across diverse tasks. Instruction-tuned models	148
	like Thought2Text (Mishra et al., 2025) further	149
	enable open-vocabulary text generation from neu-	150
	ral signals. Beyond text, LLMs serve as semantic	151
	intermediaries for cross-modal generation, guid-	152
	ing diffusion-based image synthesis (Liu et al.,	153
	2024a; Xie et al., 2025d) and 3D object reconstruc-	154
	tion (Deng et al., 2025) from brain activity. These	155
	developments demonstrate the potential of lever-	156
	aging language model architectures and training	157
	paradigms to advance neural signal analysis and	158
	brain-computer interface applications.	159
	<b>LLM-based Emotion Analysis</b> Large Language	160
	Models (LLMs) have shown strong capabilities	161
	in emotion analysis through their language under-	162
	standing and reasoning abilities. Early LLM-based	163
	approaches focused on textual emotion recognition,	164
	where text encoders extract embeddings projected	165
	into a unified space for emotion classification (Lei	166
	et al., 2023; Ma et al., 2025b,a). These methods	167
	achieved promising zero-shot and few-shot perfor-	168
	mance using prompting strategies such as instruc-	169
	tion following, in-context learning, and chain-of-	170
	thought reasoning (Wei et al., 2022). However,	171
	their reliance on text alone limits their ability to cap-	172
	ture critical affective cues from visual and acoustic	173
	modalities, such as facial micro-expressions (Xie	174
	et al., 2025a; Lin et al., 2024; Xie et al., 2025c)	175
	and prosodic variations (Feng et al., 2025a,b). This	176
	has motivated the development of multimodal large	177
	language models (MLLMs) that integrate heteroge-	178
	neous signals through specialized connector mod-	179
	ules for end-to-end multimodal emotion recogni-	180
	tion and reasoning (Yang et al., 2025b).	181
	<b>Multimodal Large Language Models</b> The	182
	emergence of Multimodal Large Language Mod-	183
	els has marked a pivotal shift in AI, introducing	184
	capabilities that span both language and vision.	185

Task	Example	Train	Eval
Individual Emotion Description (IED)	Describe the emotion of <EEG>.	✓	✓
Emotion Pairwise Comparison (EPC)	Is the individual with <EEG 1> happier than the individual with <EEG 2>?	✓	✓
Emotion Superlative Selection (ESS)	Given three EEG segments: a) <EEG 1>, b) <EEG 2>, c) <EEG 3>, select the saddest.	✓	✓
Emotion Individual Matching (EIM)	Match three EEG segments to: 1) someone waiting for a bus, 2) a parent holding their newborn, 3) a person startled by loud noise.	✓	✓
Emotion Scenario Reasoning (ESR)	<i>T1</i> : Describe three EEG segments. <i>T2</i> : [Scenario description] Select the most appropriate segment with justification.	✗	✓

Table 1: Overview of five task types for EEG-based emotion understanding.

The field’s evolution began with leveraging LLMs as coordination systems for task-specific applications (Shen et al., 2023; Yang et al., 2023; Xie et al., 2025b), then progressed toward lightweight adaptation techniques (Hu et al., 2022) and instruction-driven alignment strategies (Liu et al., 2024b) that connect visual and textual semantics. In our work, we develop the first large language model capable of bridging EEG signals and emotion analysis, addressing the challenge of automated emotion recognition from neural data.

### 3 Method

#### 3.1 Data Construction

We construct our training and evaluation data based on the SEED-VII dataset (Jiang et al., 2024a), a multimodal emotion recognition benchmark containing EEG signals from 20 subjects across seven emotion categories: happiness, surprise, neutrality, disgust, fear, sadness, and anger. The original dataset provides 62-channel EEG recordings sampled at 1000 Hz.

We first apply a 0.1–70 Hz bandpass filter and a 50 Hz notch filter to these raw signals. In alignment with EEGPT (Wang et al., 2024), we utilize 58 designated electrodes. The EEG signals are resampled to 256 Hz and partitioned into 10-second segments. Subsequently, they are scaled to mV, re-referenced using a global average reference, and extracted into 4-second windows using a random cropping strategy. Following this preprocessing pipeline, we obtain a total of 24,444 training segments and 2,761 test segments organized by emotion category.

To enable the LLM to perform interpretive emotional analysis from EEG signals, we design a diverse set of question-answer templates spanning five task types with increasing complexity, as sum-

marized in Table 1.

**Individual Emotion Description (IED):** Given a single EEG segment, the model must describe the underlying emotional state (e.g., "Describe the emotion of <EEG>."). The target response provides both an unstructured interpretation (e.g., "The EEG signals indicate a state of depression, often associated with feelings of loss or disappointment") and an explicit emotion label.

**Emotion Pairwise Comparison (EPC):** The model receives two EEG segments and must determine comparative emotional intensity (e.g., "Is the individual with EEG signals <EEG 1> happier than the individual with <EEG 2>?"). This requires the model to first describe both emotional states before drawing a comparative conclusion.

**Emotion Superlative Selection (ESS):** Given three EEG segments, the model identifies which exhibits the strongest or weakest manifestation of a target emotion (e.g., "Given three segments of EEG signals: a) <EEG 1>, b) <EEG 2>, c) <EEG 3>, select the saddest."), requiring multi-sample reasoning and ranking.

**Emotion Individual Matching (EIM):** The model matches three EEG segments to scenario descriptions of individuals experiencing different emotions (e.g., matching segments to "someone waiting for a bus on an ordinary afternoon," "a new parent holding their sleeping newborn," and "a person startled by an abrupt, loud noise"), evaluating the model’s ability to ground physiological signals in real-world emotional contexts.

**Emotion Scenario Reasoning (ESR):** This two-turn task first requires the model to describe three EEG segments, then presents a nuanced scenario (e.g., "At a family gathering, your aunt sets the table with one fewer plate than usual. No one comments, but your cousin quietly takes the extra chair away...") and asks the model to select the most appropriate segment with justification. This task evaluates complex reasoning by requiring implicit emotion inference from contextual cues.

We generate 10,000 training samples using the first four task types (IED, EPC, ESS, EIM) derived from the processed training segments, ensuring balanced representation across emotion categories. For evaluation, we utilize the 2,761 held-out test segments to construct 2,761 IED samples, 500 multi-task reasoning (EPC, ESS, and EIM) samples, and 167 ESR samples. These are held out from training to assess the model’s generalization to complex emotional reasoning scenarios. All

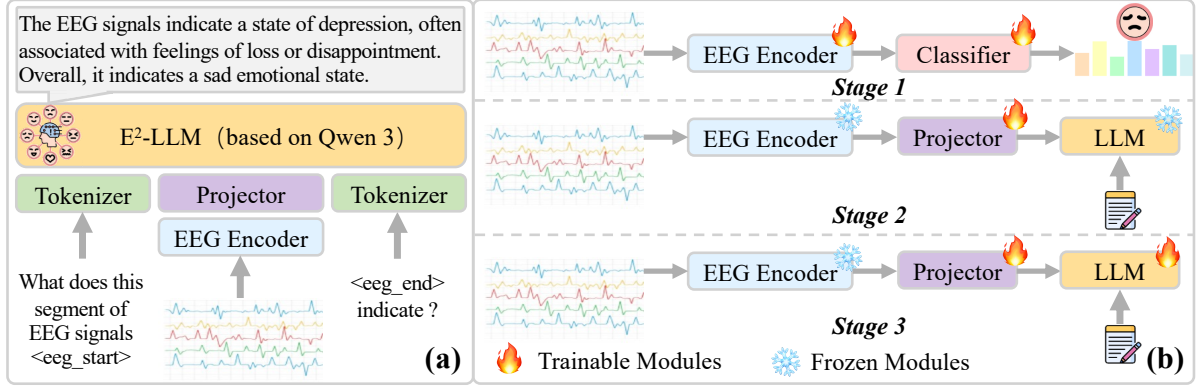


Figure 1: Overview of the E<sup>2</sup>-LLM framework and training pipeline. (a) E<sup>2</sup>-LLM framework: EEG signals are processed by an EEG Encoder, whose representations are mapped to the embedding space of a Qwen3-based LLM via a Projector. Special tokens <eeg\_start> and <eeg\_end> demarcate EEG segments within the input sequence, enabling the LLM to generate interpretive emotional analysis. (b) Multi-Stage Training Strategy: Stage 1 trains the EEG Encoder with a classification objective for emotion recognition. Stage 2 freezes the encoder and trains the Projector to align EEG representations with LLM embeddings. Stage 3 jointly fine-tunes the Projector and LLM for generating natural language emotional analysis reports.

question-answer pairs follow a chain-of-thought format, first describing individual emotional states before providing final conclusions, which encourages interpretable reasoning during inference.

## 3.2 E<sup>2</sup>-LLM

### 3.2.1 Overview

Figure 1 illustrates the overall architecture of E<sup>2</sup>-LLM, which comprises three core components: an EEG Encoder based on EEGPT (Wang et al., 2024), a Projector, and a Qwen3 (Yang et al., 2025a)-based Large Language Model (LLM). The framework processes EEG signals to generate interpretive emotional analysis through a unified multimodal architecture.

**EEG Signal Processing** Given an input EEG signal  $\mathbf{X} \in \mathbb{R}^{M \times T}$  with  $M$  channels and  $T$  time points, we first segment it into non-overlapping patches  $\mathbf{p}_{i,j}$  along both spatial and temporal dimensions:

$$\mathbf{p}_{i,j} = \mathbf{X}_{i,(j-1)d:j,d}, \quad (1)$$

$$i \in \{1, \dots, M\}, j \in \{1, \dots, N\}$$

where  $d$  denotes the temporal patch length and  $N = T/d$  represents the number of temporal patches. Each patch is then embedded through a local spatio-temporal embedding layer:

$$\mathbf{e}_{i,j} = \mathbf{W}_p^\top \mathbf{p}_{i,j} + \mathbf{b}_p + \varsigma_i \quad (2)$$

where  $\mathbf{W}_p \in \mathbb{R}^{d \times d_e}$  and  $\mathbf{b}_p \in \mathbb{R}^{d_e}$  are learnable parameters, and  $\varsigma_i$  denotes the channel-specific embedding retrieved from a learnable codebook.

**Hierarchical Encoding** The EEG Encoder adopts a hierarchical transformer architecture that processes spatial and temporal information separately. For each time step  $j$ , the encoder aggregates spatial information across channels:

$$\mathbf{h}_j = \text{ENC}(\mathbf{e}_{i,j}_{i=1}^M) \quad (3)$$

This design reduces computational complexity from  $\mathcal{O}((M \times N)^2)$  to  $\mathcal{O}(M^2 \times N)$  while enhancing flexibility for varying electrode configurations.

**Multimodal Alignment** The Projector maps EEG representations into the embedding space of the language model. We employ a two-layer MLP with GELU activation:

$$\mathbf{z} = \mathbf{W}_2 \cdot \text{GELU}(\mathbf{W}_1 \cdot \mathbf{h} + \mathbf{b}_1) + \mathbf{b}_2 \quad (4)$$

where  $\mathbf{h}$  denotes the encoder output and  $\mathbf{z} \in \mathbb{R}^{d_{\text{LLM}}}$  represents the projected embedding aligned with the LLM’s token space.

**Input Formulation** We introduce special tokens <eeg\_start> and <eeg\_end> to demarcate EEG segments within the input sequence. The final input to LLM is constructed as:

$$\mathbf{S} = [\mathbf{w}_1, \dots, \mathbf{w}_k, \langle \text{eeg\_start} \rangle, \mathbf{z}_1, \dots, \mathbf{z}_N, \langle \text{eeg\_end} \rangle, \mathbf{w}_{k+1}, \dots] \quad (5)$$

where  $\mathbf{w}_i$  are text token embeddings and  $\mathbf{z}_j$  are projected EEG embeddings, enabling the LLM to generate natural language emotional analysis reports conditioned on both textual instructions and EEG signals.

### 3.2.2 Multi-stage Training Strategy

As illustrated in Figure 1(b), we adopt a three-stage training pipeline to progressively bridge the modality gap between EEG signals and natural language. This curriculum-based approach ensures stable optimization and effective cross-modal alignment.

**Stage 1: EEG Encoder Training** In the first stage, we fine-tune the pretrained EEG Encoder with a classification objective to learn discriminative emotion representations. Given the encoded features  $\mathbf{h} = \text{ENC}(\mathbf{X})$ , a classification head predicts the emotion category:

$$\hat{y} = \text{Softmax}(\mathbf{W}_c \cdot \mathbf{h} + \mathbf{b}_c) \quad (6)$$

where  $\mathbf{W}_c \in \mathbb{R}^{C \times d_e}$  and  $\mathbf{b}_c \in \mathbb{R}^C$  are learnable parameters, and  $C$  denotes the number of emotion categories. The encoder is optimized using cross-entropy loss:

$$\mathcal{L}_{\text{cls}} = - \sum_{c=1}^C y_c \log(\hat{y}_c) \quad (7)$$

where  $y_c$  represents the ground-truth label. During this stage, both the Projector and LLM modules remain frozen, allowing the encoder to focus on extracting emotion-relevant features from raw EEG signals.

**Stage 2: Cross-modal Alignment** In the second stage, we freeze the EEG Encoder and train the Projector to align EEG representations with the LLM’s embedding space. The model is trained using the autoregressive language modeling objective:

$$\mathcal{L}_{\text{LM}} = - \sum_{l=1}^L \log P_{\theta}(w_l | \mathbf{S}_{<l}) \quad (8)$$

where  $L$  is the length of the target response,  $w_l$  denotes the  $l$ -th token, and  $\mathbf{S}_{<l}$  represents all preceding tokens including both textual instructions and projected EEG embeddings. This stage establishes a semantic bridge between the two modalities while preserving the learned EEG representations.

**Stage 3: End-to-end Fine-tuning** In the final stage, we jointly fine-tune the Projector and LLM while keeping the EEG Encoder frozen. The model continues to be optimized with the same language modeling objective  $\mathcal{L}_{\text{LM}}$ , but now updates both the Projector and LLM parameters. To enhance parameter efficiency, we employ Low-Rank Adaptation (LoRA) (Hu et al., 2022; Zhu et al., 2024) for the

LLM, which significantly reduces the number of trainable parameters while maintaining model expressiveness. This stage enables the model to generate fluent and accurate natural language emotional analysis reports conditioned on EEG inputs.

## 4 Experiments

### 4.1 Experimental Setup

Our architecture incorporates the 10M-parameter EEG encoder from EEGPT (Wang et al., 2024) and adopt Qwen3 (Yang et al., 2025a) as the backbone language model. We train and evaluate three variants of the proposed framework—E<sup>2</sup>-LLM<sub>4B</sub>, E<sup>2</sup>-LLM<sub>8B</sub>, and E<sup>2</sup>-LLM<sub>14B</sub>—employing Qwen3-4B, Qwen3-8B, and Qwen3-14B respectively.

**Training Details** Encoder fine-tuning was conducted on the 24,444 training segments for 50 epochs using the AdamW optimizer (Loshchilov and Hutter, 2017) with a weight decay of 0.01, a learning rate of  $10^{-4}$ , batch size of 64, following a one-cycle learning rate schedule. For cross-modal alignment, the Projector module was trained on 10,000 training question-answer samples using Adam (Kingma, 2014) with no weight decay, a learning rate of  $2 \times 10^{-4}$ , and a batch size of 16 under a cosine annealing schedule.

For end-to-end fine-tuning, both the Projector and the LoRA (Hu et al., 2022) parameters of the LLM were optimized using 3,000 training samples. We maintained a batch size of 16 and a learning rate of  $2 \times 10^{-4}$  for both components, utilizing Adam under a cosine annealing schedule without weight decay. The LoRA configuration involved a scaling factor of 256, a rank of 128, and a dropout rate of 0.05, targeting the Query Projection Layer ( $W_Q$ ) and the Key Projection Layer ( $W_K$ ) of the language model.

Regarding computational costs, encoder fine-tuning took less than 7 hours and required less than 34 GB of GPU VRAM. The combined cross-modal alignment and end-to-end fine-tuning stages required 2.3, 3.2, and 5.4 hours for the 4B, 8B, and 14B variants, respectively. Experiments for the 4B model were conducted on a single NVIDIA RTX 6000 Ada Generation GPU, while the 8B and 14B models were trained using two such GPUs.

**Evaluation Protocol** To rigorously evaluate the performance of E<sup>2</sup>-LLM, we utilize the test QA pairs derived from the held-out test segments. The evaluation framework spans three primary dimen-

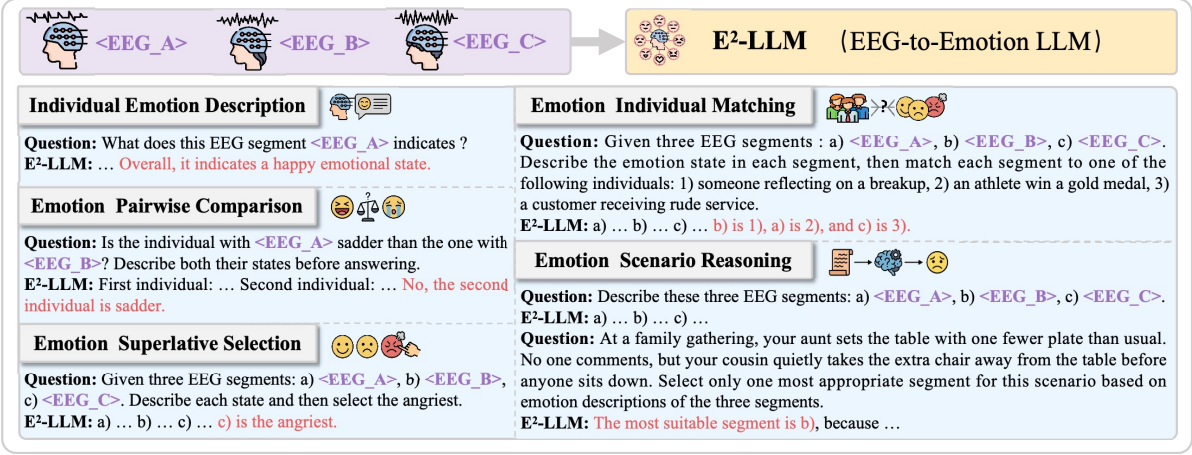


Figure 2: Illustrative examples from the proposed E<sup>2</sup>-LLM across five distinct tasks. The qualitative results demonstrate that E<sup>2</sup>-LLM can generate reasonable outputs.

	Random	NeuroLM	E <sup>2</sup> -LLM <sub>4B</sub>	E <sup>2</sup> -LLM <sub>8B</sub>	E <sup>2</sup> -LLM <sub>14B</sub>
"Happy"	14.29	11.71	33.02	46.42	<b>51.40</b>
"Surprise"	14.29	35.62	71.82	65.34	<b>76.20</b>
"Neutral"	14.29	18.24	49.01	<b>80.24</b>	70.75
"Disgust"	14.29	14.30	8.24	<b>45.88</b>	20.88
"Fear"	14.29	40.90	<b>77.48</b>	75.00	76.35
"Sad"	14.29	35.02	<b>75.37</b>	46.56	53.65
"Anger"	14.29	21.97	67.46	72.68	<b>78.15</b>
<b>Balanced Accuracy</b>	14.29	25.39	54.63	<b>61.73</b>	61.05
<b>Cohen's Kappa</b>	0.00	13.65	49.69	54.72	<b>55.14</b>
<b>Weighted F1-score</b>	14.29	25.83	54.56	<b>60.91</b>	60.16

Table 2: Performance (%) comparison on the Individual Emotion Description (IED) task. Results demonstrate the effectiveness of E<sup>2</sup>-LLM variants against baselines across seven emotion categories. Cohen’s kappa measures prediction reliability by accounting for chance agreement. Best results are **bolded**.

sions: 1) For the basic emotion prediction task (IED), we report balanced accuracy, Cohen’s kappa, and weighted F1-score across the entire test set to provide a robust measure of classification performance. 2) We measure the accuracy on multi-task reasoning questions (EPC, ESS, and EIM) to assess the model’s multi-faceted cognitive reasoning capabilities. 3) The model’s generalization ability is quantified via accuracy over a specialized scenario-based reasoning task (ESR) that remains unseen during training.

Visual illustrations of the five tasks mentioned above and the corresponding model responses are provided in Figure 2.

## 4.2 Results

**Emotion Prediction** Table 2 summarizes the performance of the E<sup>2</sup>-LLM variants and baseline methods on the individual emotion description (IED) task. Crucially, we benchmark against NeuroLM (Jiang et al., 2024b), a recent paradigm

that similarly couples EEG encoder and LLM. Following its protocol, we fine-tuned NeuroLM in different parameter scales on SEED-VII (Jiang et al., 2024a) and report the best-performing version. However, it achieves a balanced accuracy of only 25.39%. Further analysis reveals that NeuroLM suffers from overfitting to the rigid instruction formats used during classification training, rendering it unable to generalize to varied instruction forms.

In contrast, E<sup>2</sup>-LLM demonstrates superior performance across all metrics, with the benefits of scaling the LLM backbone becoming immediately apparent. Specifically, the 8B model offers a substantial improvement over the 4B variant, attaining the highest balanced accuracy (61.73%) and weighted F1-score (60.91%). Notably, while the 14B variant shows a marginal decrease in balanced accuracy, it secures the superior Cohen’s kappa value (55.14%). Given that Cohen’s kappa is a robust measure of inter-rater agreement that ex-

	EPC	ESS	EIM	ESR*	Avg.
Random	33.33	33.33	16.67	33.33	29.17
E <sup>2</sup> -LLM <sub>4B</sub>	67.48	72.94	<b>77.84</b>	34.13	63.10
E <sup>2</sup> -LLM <sub>8B</sub>	68.71	76.47	71.86	41.92	64.74
E <sup>2</sup> -LLM <sub>14B</sub>	<b>72.39</b>	<b>79.41</b>	73.65	<b>53.89</b>	<b>69.84</b>

Table 3: Accuracy (%) results on the EPC, ESS, EIM, and ESR tasks. The ESR task, marked with \*, evaluates zero-shot generalization to unseen complex scenarios. Best results are **bolded**.

462 plicitly corrects for chance, this result underscores  
463 that the 14B model provides the most consistent  
464 and reliable prediction quality. These observations  
465 suggest that scaling the LLM backbone enhances  
466 the intrinsic reliability of EEG-based emotion de-  
467 coding, establishing larger variants as more robust  
468 foundations for classification tasks.

469 **Multi-task Reasoning** Table 3 presents the eval-  
470 uation results for multi-task reasoning capabili-  
471 ties, encompassing Emotion Pairwise Comparison  
472 (EPC), Emotion Superlative Selection (ESS), and  
473 Emotion Individual Matching (EIM). All E<sup>2</sup>-LLM  
474 models drastically surpass the corresponding ran-  
475 dom baselines across these complex cognitive tasks,  
476 decisively validating the framework’s ability to ex-  
477 tract and utilize EEG information for sophisticated  
478 relational and contextual reasoning. In general,  
479 model scaling correlates positively with enhanced  
480 reasoning capability, with E<sup>2</sup>-LLM<sub>14B</sub> achiev-  
481 ing the highest overall average accuracy. This out-  
482 come supports the hypothesis that increased LLM ca-  
483 pacity provides a superior foundation for handling  
484 intricate, diverse cognitive demands.

485 However, a deviation is observed: E<sup>2</sup>-LLM<sub>4B</sub>  
486 achieves the peak performance on the EIM task  
487 (77.84%), surpassing its larger 8B and 14B coun-  
488 terparts. This result suggests that a smaller model  
489 size may occasionally find an efficient, task-specific  
490 optimal mapping for tasks requiring focused se-  
491 mantic matching, which may be slightly diluted by  
492 the broader reasoning capacities optimized in the  
493 larger models.

494 **Scenario Reasoning** The Emotion Scenario Rea-  
495 soning (ESR) task serves as a critical test of gener-  
496 alization, as its complex two-turn format and con-  
497 textual inference requirements were unseen during  
498 training. As shown in Table 3, this task reveals a  
499 distinct phenomenon linked to model scale. The  
500 4B model cannot generalize well, yielding an accu-  
501 racy of 34.13%, which is just slightly higher than

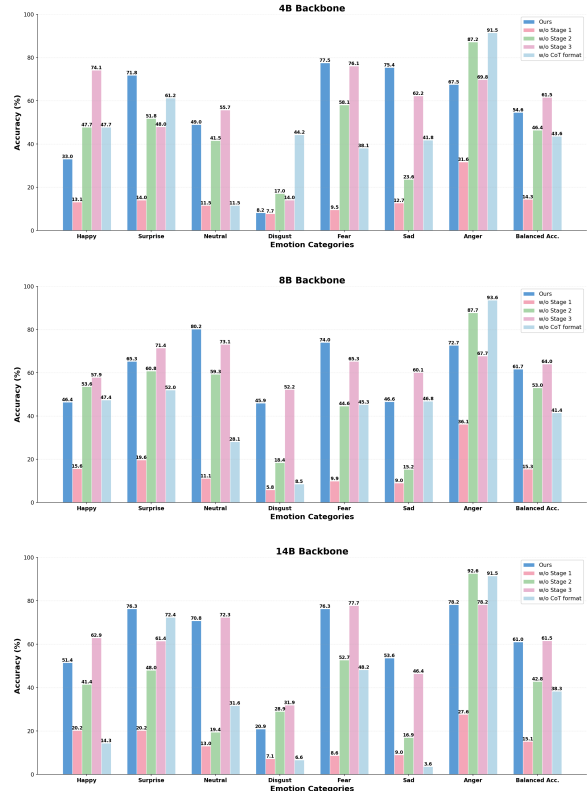


Figure 3: Ablation analysis on the IED task across different model scales.

502 the random baseline. Qualitative inspection indi-  
503 cates that the 4B model sometimes struggles to  
504 adhere to the novel instruction format, often fail-  
505 ing to generate valid responses—a limitation typi-  
506 cal of smaller models that overfit to fixed training  
507 templates, which is even much more severe in Neu-  
508 roLM (Jiang et al., 2024b). Conversely, increas-  
509 ing the model size to 8B and 14B unlocks signifi-  
510 cant zero-shot generalization capabilities, with per-  
511 formance jumping to 41.92% and 53.89% respec-  
512 tively. This finding suggests that while smaller models  
513 can master specific EEG-to-text mappings, a larger  
514 language backbone is indispensable for flexibly ap-  
515 plying these representations to unseen, high-level  
516 reasoning scenarios.

### 4.3 Ablation Studies 517

518 To validate the contribution of each training stage  
519 and the chain-of-thought question-answer format,  
520 we conduct ablation studies on the five tasks (IED,  
521 EPC, ESS, EIM, and ESR), with results summa-  
522 rized in Figure 3 and Table 4. 522

523 We analyze four variants: 1) "w/o Stage 1": Us-  
524 ing the pretrained EEGPT encoder (Wang et al.,  
525 2024) directly instead of the training it with a clas-  
526 sification objective for emotion recognition; 2) "w/o

	EPC	ESS	EIM	ESR	Avg.
<b>E<sup>2</sup>-LLM<sub>4B</sub></b>	<b>67.5</b>	<b>72.9</b>	<b>77.8</b>	34.1	<b>63.1</b>
w/o Stage 1	32.5	28.2	14.4	7.2	20.6
w/o Stage 2	66.3	72.9	66.5	37.1	60.7
w/o Stage 3	65.6	70.0	68.9	15.0	54.9
w/o CoT format	55.8	65.3	62.9	<b>38.3</b>	55.6
<b>E<sup>2</sup>-LLM<sub>8B</sub></b>	68.7	76.5	<b>71.9</b>	<b>41.9</b>	<b>64.7</b>
w/o Stage 1	42.3	28.8	17.4	28.7	29.3
w/o Stage 2	68.7	<b>77.1</b>	69.5	37.7	63.3
w/o Stage 3	<b>74.2</b>	73.5	59.9	25.7	58.3
w/o CoT format	53.4	31.2	21.0	40.1	36.4
<b>E<sup>2</sup>-LLM<sub>14B</sub></b>	<b>72.4</b>	<b>79.4</b>	73.7	<b>53.9</b>	<b>69.8</b>
w/o Stage 1	39.9	30.6	21.0	31.7	30.8
w/o Stage 2	67.5	75.9	71.3	40.1	63.7
w/o Stage 3	67.5	65.9	64.7	38.9	59.3
w/o CoT format	58.9	69.4	<b>76.6</b>	50.9	64.0

Table 4: Ablation analysis on the EPC, ESS, EIM, and ESR tasks across different model scales. We **bold** the best for each backbone.

Stage 2": Skipping the separate alignment training for the Projector module and proceeding directly to end-to-end fine-tuning; 3) "w/o Stage 3": Using the model after the alignment stage without end-to-end instruction tuning; and 4) "w/o CoT format": Training and evaluating the model without the prompts and reference answers that guide the LLM to analyze emotion states before providing conclusion.

**Impact of Multi-stage Training** The results demonstrate that the proposed multi-stage training strategy is critical for effective emotion recognition and reasoning. First, removing the emotion-discriminative training for the EEG encoder ("w/o Stage 1") leads to a catastrophic performance drop across all tasks and backbone sizes, with balanced accuracy falling to near-chance levels. This indicates that the LLM cannot compensate for a "blind" encoder, which means the EEG encoder must learn robust emotion-specific representations before being aligned with the LLM. Second, the cross-modal alignment is equally vital—the "w/o Stage 2" variant suffers a significant performance degradation (e.g., dropping from 61.0% to 42.8% for the 14B model on the IED task). This confirms that the Projector requires a dedicated alignment phase to map physiological signals into the semantic space of the LLM.

An interesting observation from the IED task is that the model without instruction tuning ("w/o Stage 3") achieves classification accuracy comparable to, or marginally higher than, the full model

(e.g., 61.5% vs. 61.0% for the 14B backbone). This suggests that for simple classification tasks (IED), the alignment learned in Stage 2 is sufficient. However, as shown in Table 4, the "w/o Stage 3" model fails significantly on more complex reasoning tasks (EPC, ESS, EIM, and ESR), lacking the instruction-following capabilities required to compare emotions or associate emotions with context. Stage 3 is therefore essential for generalizing from simple recognition to interpretable analysis, even if it introduces a slight trade-off in pure classification metrics due to the increased complexity of the generation objective.

**Efficacy of Chain-of-Thought** The removal of the chain-of-thought format in both training and evaluation question-answer samples ("w/o CoT format") results in a sharp decline in overall accuracy across all scales (e.g., an average drop of almost 10% for the 14B model across five tasks). This validates our hypothesis that decomposing the task—forcing the model to analyze and describe the emotion state of EEG before predicting the label or selecting the answer—serves as a crucial regularizer, enabling the LLM to ground its predictions in the observed signal dynamics rather than hallucinating answers.

## 5 Conclusion

We present E<sup>2</sup>-LLM, the first multimodal large language model framework that bridges EEG signals with interpretable emotion analysis. By integrating a pretrained EEG encoder with Qwen-based LLMs through learnable projections and a multi-stage training pipeline, E<sup>2</sup>-LLM achieves strong emotion recognition performance while providing human-interpretable explanations of affective states. Our evaluation demonstrates that model scaling enhances both classification accuracy and zero-shot generalization to complex reasoning scenarios. This work establishes a new paradigm combining physiological signals with LLM reasoning capabilities, showing that neural dynamics can be effectively translated into semantic emotional understanding. Future directions include incorporating additional modalities, improving cross-dataset generalization ability, and addressing the inherent subjectivity in emotion labeling for real-world deployment.

## 606 Limitation

607 Our study has several limitations that should be  
608 acknowledged. First, experiments are conducted  
609 solely on the SEED-VII dataset with only 20 sub-  
610 jects, which may limit generalizability across di-  
611 verse populations and neurophysiological varia-  
612 tions. Second, the larger model variants (8B and  
613 14B) require substantial computational resources,  
614 potentially limiting accessibility for researchers  
615 with constrained budgets. Third, we observe incon-  
616 sistent scaling behaviors across tasks—notably, the  
617 smaller 4B model outperforms larger variants on  
618 the Emotion Individual Matching task, suggesting  
619 non-monotonic relationships between model size  
620 and performance. Fourth, even our best-performing  
621 model achieves only 53.89% accuracy on zero-shot  
622 scenario reasoning, indicating significant room for  
623 improvement in generalizing to unseen complex  
624 emotional contexts. Finally, while E<sup>2</sup>-LLM gener-  
625 ates natural language explanations, we lack system-  
626 atic human evaluation of the semantic quality and  
627 clinical relevance of these interpretations, which  
628 remains an important direction for establishing rig-  
629 orous assessment metrics.

## 630 Ethical Concern

631 This research involves analysis of EEG signals  
632 for emotion recognition, which raises important  
633 ethical considerations regarding privacy and po-  
634 tential misuse. EEG data contains sensitive neu-  
635 rophysiological information that could reveal pri-  
636 vate mental states beyond intended emotion cat-  
637 egories. While our experiments use the publicly  
638 available SEED-VII dataset with appropriate in-  
639 formed consent, we acknowledge that automated  
640 emotion recognition systems could be misused for  
641 unauthorized surveillance, workplace monitoring,  
642 or discriminatory decision-making without explicit  
643 user consent. Furthermore, the model’s interpretive  
644 capabilities could potentially generate seemingly  
645 authoritative but inaccurate psychological assess-  
646 ments. We strongly advocate that deployment of  
647 EEG-based emotion recognition must be governed  
648 by strict ethical guidelines, including explicit in-  
649 formed consent, transparent data usage disclosure,  
650 robust privacy protections, and safeguards against  
651 coercive applications. Future work should prior-  
652 itize developing technical safeguards and regula-  
653 tory frameworks to ensure neural signal analysis  
654 respects human dignity and individual rights.

## References 655

- 656 Naseem Babu, Jimson Mathew, and AP Vinod. 2025. 656  
657 Large language models for eeg: A comprehensive sur- 657  
658 vey and taxonomy. *arXiv preprint arXiv:2506.06353*. 658
- 659 Hongyu Chen, Weiming Zeng, Chengcheng Chen, 659  
660 Luhui Cai, Fei Wang, Yuhu Shi, Lei Wang, Wei 660  
661 Zhang, Yueyang Li, Hongjie Yan, and 1 others. 2025. 661  
662 Eeg emotion copilot: Optimizing lightweight llms for 662  
663 emotional eeg interpretation with assisted medical 663  
664 record generation. *Neural Networks*, page 107848. 664
- 665 Xia Deng, Shen Chen, Jiale Zhou, and Lei Li. 2025. 665  
666 Mind2matter: Creating 3d models from eeg signals. 666  
667 *arXiv preprint arXiv:2504.11936*. 667
- 668 Tao Feng, Yifan Xie, Xun Guan, Jiyuan Song, Zhou 668  
669 Liu, Fei Ma, and Fei Yu. 2025a. Unisync: A unified 669  
670 framework for audio-visual synchronization. *arXiv 670*  
671 *preprint arXiv:2503.16357*. 671
- 672 Tao Feng, Zhiyuan Zhao, Yifan Xie, Yuqi Ye, Xi- 672  
673 angyang Luo, Xun Guan, and Yu Li. 2025b. Stft- 673  
674 codec: High-fidelity audio compression through time- 674  
675 frequency domain representation. *arXiv preprint 675*  
676 *arXiv:2503.16989*. 676
- 677 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan 677  
678 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 678  
679 Weizhu Chen, and 1 others. 2022. Lora: Low-rank 679  
680 adaptation of large language models. *ICLR*, 1(2):3. 680
- 681 Gan Huang, Zhiheng Zhao, Shaorong Zhang, Zhenxing 681  
682 Hu, Jiaming Fan, Meisong Fu, Jiale Chen, Yaqiong 682  
683 Xiao, Jun Wang, and Guo Dan. 2023. Discrepan- 683  
684 cy between inter-and intra-subject variability in 684  
685 eeg-based motor imagery brain-computer interface: 685  
686 Evidence from multiple perspectives. *Frontiers in 686*  
687 *neuroscience*, 17:1122661. 687
- 688 Wei-Bang Jiang, Xuan-Hao Liu, Wei-Long Zheng, and 688  
689 Bao-Liang Lu. 2024a. Seed-vii: A multimodal 689  
690 dataset of six basic emotions with continuous labels 690  
691 for emotion recognition. *IEEE Transactions on Af- 691*  
692 *fective Computing*. 692
- 693 Wei-Bang Jiang, Yansen Wang, Bao-Liang Lu, and 693  
694 Dongsheng Li. 2024b. Neurolm: A universal multi- 694  
695 task foundation model for bridging the gap be- 695  
696 tween language and eeg signals. *arXiv preprint 696*  
697 *arXiv:2409.00101*. 697
- 698 Weibang Jiang, Liming Zhao, and Bao-liang Lu. 2024c. 698  
699 Large brain model for learning generic representa- 699  
700 tions with tremendous eeg data in bci. In *The Twelfth 700*  
701 *International Conference on Learning Representa-* 701  
702 *tions*. 702
- 703 Diederik P Kingma. 2014. Adam: A method for stochas- 703  
704 tic optimization. *arXiv preprint arXiv:1412.6980*. 704
- 705 Shanglin Lei, Guanting Dong, Xiaoping Wang, Keheng 705  
706 Wang, Runqi Qiao, and Sirui Wang. 2023. Instruc- 706  
707 terc: Reforming emotion recognition in conversation 707  
708 with multi-task retrieval-augmented large language 708  
709 models. *arXiv preprint arXiv:2309.11911*. 709

