

---

# EEG2TEXT: Open Vocabulary EEG-to-Text Decoding with EEG Pre-Training and Multi-View Transformer

---

Hanwen Liu<sup>1</sup> Daniel Hajjaligol<sup>1</sup> Benny Antony<sup>2</sup> Aiguo Han<sup>3</sup> Xuan Wang<sup>1</sup>

## Abstract

Deciphering the intricacies of the human brain has captivated curiosity for centuries. Recent strides in Brain-Computer Interface (BCI) technology, particularly using motor imagery, have restored motor functions such as reaching, grasping, and walking in paralyzed individuals. However, unraveling natural language from brain signals remains a formidable challenge. Electroencephalography (EEG) is a non-invasive technique used to record electrical activity in the brain by placing electrodes on the scalp. Previous studies of EEG-to-text decoding have achieved high accuracy on small closed vocabularies, but still fall short of high accuracy when dealing with large open vocabularies. We propose a novel method, EEG2TEXT, to improve the accuracy of open vocabulary EEG-to-text decoding. Specifically, EEG2TEXT leverages EEG pre-training to enhance the learning of semantics from EEG signals and proposes a multi-view transformer to model the EEG signal processing by different spatial regions of the brain. Experiments show that EEG2TEXT has superior performance, outperforming the state-of-the-art baseline methods by a large margin of up to 5% in absolute BLEU and ROUGE scores. EEG2TEXT shows great potential for a high-performance open-vocabulary brain-to-text system to facilitate communication.

## 1. Introduction

Recent advances in brain-computer interface (BCI) technology have demonstrated exciting progress in restoring

---

<sup>\*</sup>Equal contribution <sup>1</sup>Department of Computer Science, Virginia Tech, VA, USA <sup>2</sup>Department of Electrical and Computer Engineering, Virginia Tech, VA, USA <sup>3</sup>Department of Biomedical Engineering and Mechanics, Virginia Tech, VA, USA. Correspondence to: Hanwen Liu <liuhwen@vt.edu>, Daniel Hajjaligol <danielhajjaligol@vt.edu>, Benny Antony <bennyantony@vt.edu>, Aiguo Han <aiguohan@vt.edu>, Xuan Wang <xuanw@vt.edu>.

the capabilities of patients with paralysis, such as reaching (Hochberg et al., 2012), grasping (Aflalo et al., 2015; Boulton et al., 2016), and walking (Lorach et al., 2023). The heart of BCI is its ability to accurately decode complex brain signals. Despite the advances in decoding brain signals related to motion, decoding brain signals related to speech remains a formidable challenge. Previous research translating speech-related brain signals to text (brain-to-text) primarily relies on electrocorticography (ECoG), an invasive electrophysiological monitoring method that uses electrodes placed directly on the exposed brain surface to record activity from the cerebral cortex. ECoG offers higher temporal and spatial resolution than traditional noninvasive scalp electroencephalography (EEG), with a significantly better signal-to-noise ratio. However, the invasive nature of ECoG is undesirable for BCI applications, and it is highly desirable to develop brain-to-text decoding methods using noninvasive EEG signals, although EEG signals are significantly more challenging to work with than ECoG.

Previous studies of EEG-to-text decoding (Herff et al., 2015; Sun et al., 2019; Anumanchipalli et al., 2019; Makin et al., 2020; Panachakel & Ramakrishnan, 2021; Moses et al., 2021; Nieto et al., 2022) have achieved high accuracy on small closed vocabularies, but still fall short of high accuracy when dealing with large open vocabularies. These approaches primarily target high accuracy (> 90%) but are often confined to small closed vocabularies and struggle to decode semantically similar words beyond training sets. Recent studies broaden the scope from closed to open-vocabulary EEG-to-text decoding (Wang & Ji, 2021; Willett et al., 2023; Tang et al., 2023; Duan et al., 2023), drastically expanding the vocabulary size by over 100-fold, from several hundred to tens of thousands of words. Notably, two of these studies (Wang & Ji, 2021; Duan et al., 2023) leverage a pre-trained large language model BART (Lewis et al., 2019), and represent the state-of-the-art for open vocabulary brain-to-text decoding. However, these studies are in their nascent stages and are challenged by their limited accuracy.

To improve the accuracy of EEG-to-text decoding with open vocabularies, we propose a novel EEG-to-text decoding method based on transformers. First, we introduce a Convolutional Neural Network (CNN) module before the base transformer model to enhance the model’s ability to han-

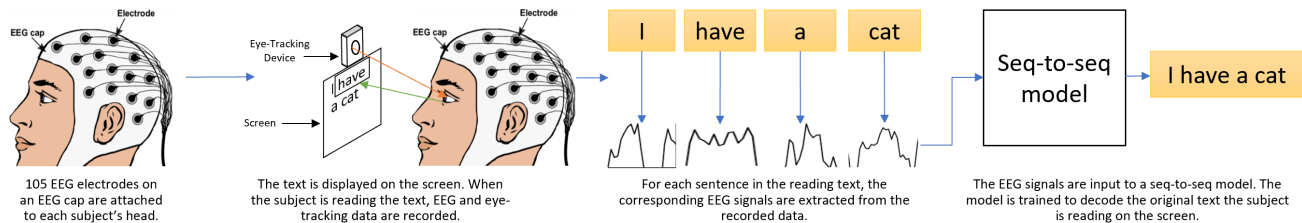


Figure 1. The overall framework of open-vocabulary EEG-to-text translation. The first sub-figure comes from (Nagel & Spüler, 2018).

dle long EEG signals. Second, we conduct pre-training of the transformer model by reconstructing randomly masked EEG signals from the input data. This pre-training step helps our transformer model better learn the semantics of EEG signals. Last, we propose a multi-view transformer architecture, where each single-view transformer is the pre-trained model from the previous step, to model the EEG signal processing by different spatial regions of the brain. Experiments show that EEG2TEXT has superior performance, outperforming the state-of-the-art baseline methods by a large margin of up to 5% in absolute BLEU and ROUGE scores. EEG2TEXT shows great potential for a high-performance open-vocabulary brain-to-text system to facilitate communication. We will open-source our code and dataset to facilitate future studies of EEG-to-text translation.

## 2. Task Definition

Our task involves decoding corresponding text from EEG signals (Figure 1). The data acquisition process involves 1) attaching an EEG cap to each subject's head, 2) displaying the text (reading materials) on a screen, and 3) recording the EEG and eye-tracking (for verification and calibration of the EEG signals) data while the subject is reading the text. The EEG signals are further extracted from the recorded data and fed as input to a decoding model to predict the original text the subject was reading on the screen.

Formally, this task can be formulated as a sequence-to-sequence machine translation task as follows:

$$P(Y|X) = \arg \max_Y \prod_{t=1}^{T'} P(y_t | y_{<t}, X) \quad (1)$$

where  $T'$  represents the length of the target sentence  $Y$ ;  $y_t$  represents the word or token at position  $t$  in the target sentence  $Y$ ;  $y_{<t}$  represents the words or tokens preceding position  $t$  in the target sentence  $Y$ ;  $X$  represents the input EEG data; and  $P(y_t | y_{<t}, X)$  is the conditional probability of generating word  $y_t$  given the previous words  $y_{<t}$  and the input EEG data  $X$ . Our goal is to maximize the probability  $P(Y|X)$  of generating the target sentence given the input EEG data.

## 3. Methodology

### 3.1. Baseline Model

Our baseline model (Wang & Ji, 2021) takes the word-level EEG features as the input to a transformer model followed by a pre-trained BART model for text decoding. The raw EEG signals are typically stored as a two-dimensional array with one dimension for time and the other for channels (the number of electrodes used to collect EEG signals). Each value in this two-dimensional array corresponds to the signal strength collected at the corresponding time for the corresponding channel. In the baseline model, the word-level EEG features are extracted from eight independent frequency bands from the raw EEG signals. The above eight word-level EEG features are simply concatenated across all the channels as input to the decoder framework.

The baseline model faces the following challenges: 1) the reliance on eye-tracking calibration for word-level EEG feature extraction introduces error propagation and lacks generalizability to scenarios such as inner speech decoding (Martin et al., 2018; Nalborczyk et al., 2020), 2) there is room for improvement in EEG representation learning through self-supervised pre-training, and 3) the lack of spatial resolution modeling ignores the varying importance of different brain regions in language processing. To overcome these challenges, we propose a novel framework, EEG2TEXT, that achieves superior performance for open-vocabulary EEG-to-text translation.

### 3.2. Convolutional Transformer for Sentence-Level EEG Encoding

Instead of using the word-level EEG features crafted based on the eye-tracking data, we directly use the sentence-level EEG signals as input to our model. Using sentence-level EEG signals offers several advantages over word-level EEG features. It provides richer information without error propagation from the eye-tracking data and exhibits better generalizability to other tasks, such as inner speech decoding, where acquiring eye-tracking data is infeasible.

However, the sentence-level EEG signals pose a challenge due to their excessive length, potentially overloading laboratory-level GPUs if directly input into the transformer

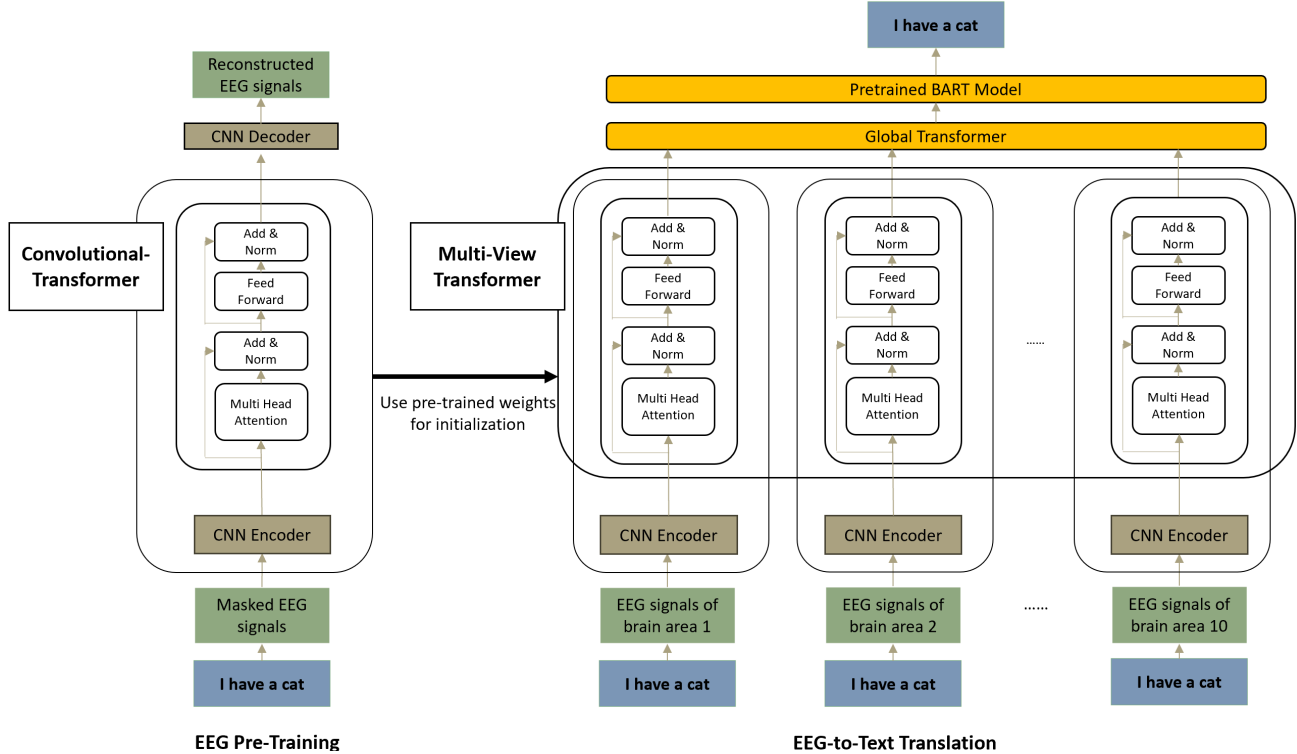


Figure 2. The overall framework of EEG2TEXT. It takes the sentence EEG signals as input and decodes the original text as output. EEG2TEXT includes major steps of 1) a base convolutional transformer model, 2) pre-training for EEG encoding, and 3) a multi-view transformer for different spatial regions of the brain.

layer. To tackle this issue, we introduce a convolutional transformer model that incorporates a CNN module for compressing raw EEG signals. Utilizing CNN-Transformer for modeling long sequences has been proven effective in previous EEG signal processing tasks (Song et al., 2022). So we choose this CNN-Transformer as the base architecture to develop our models. This CNN module comprises two convolutional layers, adept at both temporal and spatial (or channel) compression. We also compared two input formats of the sentence-level EEG signals: 1) the raw signals, and 2) the spectrogram of the signals. The spectrogram of a signal (Appendix Figure A1) is a two-dimensional image, where the x-axis represents time, the y-axis represents frequency, and the image pixel value represents the magnitude of the signal at each time-frequency pair. The sentence-level EEG signals are then input into the CNN module to obtain compressed EEG signals, which are then fed into the transformer model for subsequent feature extraction and text translation.

### 3.3. Transformer Pre-Training for an Enhanced EEG Encoding

To enhance the semantic understanding of the EEG signals, we propose a self-supervised pre-training of the convolutional transformer model for parameter initialization (Figure

2). Inspired by the masked language model pre-training strategies (Devlin et al., 2018; Joshi et al., 2019; Liu et al., 2019), we formulate our self-supervised pre-training objective as follows:

$$\theta^* = \arg \max_{\theta} \sum_{(i,j) \in \mathcal{D}} \log P(M|C; \theta), \quad (2)$$

where  $M$  represents the masked tokens;  $C$  represents the context or surrounding tokens;  $\theta^*$  represents the optimal model parameters;  $\theta$  represents the model parameters being optimized;  $\mathcal{D}$  represents the training data, where  $(i, j)$  are pairs of sentences or sentence fragments; and  $P(M|C; \theta)$  is the probability of predicting the masked tokens.

During the self-supervised pre-training stage, we add a convolutional decoder module on top of the convolutional transformer encoder to decode the input EEG signals. The input is the sentence-level EEG signals masked with different strategies and the output is the sentence-level EEG signals reconstructed by the CNN decoder. Specifically, we compared three different masking strategies for the sentence-level EEG signals as follows:

- **Masked Token Prediction** (Devlin et al., 2018): randomly masking 15% of all the tokens.

Table 1. Ten channel groups and their corresponding approximate brain areas.

Approximate Brain Areas	Corresponding Electrodes
Prefrontal Cortex	E6, E12, E5, E11, E16, E15, E20, E118, E24, E124, E26, E2, E27, E123, E3, E4, E23, E19, E22, E9, E10, E18, E28, E33, E117, E122
Premotor Cortex	CZ, E7, E106, E105, E104, E115, E114, E120, E110, E116, E121, E111, E112, E109, E13, E30
Broca’s Area	E29, E36, E35, E34
Auditory Association Area	E40, E38, E39, E43, E44, E46, E57, E58, E64
Primary Motor Cortex	E31, E80, E55, E37, E87, E93, E103, E102, E108
Primary Sensory Cortex	E54, E79, E61, E78, E62, E53, E86, E92, E98, E100, E101
Somatic Sensory Cortex	E67, E77, E71, E72, E76, E66, E84, E60, E85
Auditory Cortex	E59, E91, E97, E51
Wernicke’s Area	E41, E42, E52, E47, E45, E50
Visual Area	E65, E69, E70, E74, E75, E82, E83, E89, E90, E95, E96

- **Continuous Masked Token Prediction** (Joshi et al., 2019): randomly masking a sequence of consecutive tokens until a total of 15% of all the tokens are masked.
- **Re-Masked Token Prediction** (Liu et al., 2019): re-randomizing the masking of 15% of all the tokens for each training epoch.

It is important to highlight that our self-supervised pre-training step allows for seamless integration of EEG data from diverse tasks, including image recognition. In our experiments, we further incorporated an image EEG dataset (Gifford et al., 2022) during pre-training, aiming to showcase the model’s adaptability to EEG signals from multi-modal data and explore the potential for enhanced translation performance through the combination of EEG signals from diverse data modalities.

The goal of this pre-training step is to have the convolutional transformer learn meaningful concepts such as context, relationships, and semantics present in sentence-level EEG signals during this pre-training process. After pre-training, the parameters are saved and used as the initial parameters for the final multi-view transformer model.

### 3.4. Multi-View Transformer for Different Spatial Regions of the Brain

Another important feature of our model is the novel multi-view transformer decoder architecture we introduced that encodes different regions of the brain with a different convolutional transformer (Figure 2). The multi-view transformer model takes into account the fact that different brain regions potentially play different roles in language processing. This spatial modeling therefore can improve the model performance, but has been overlooked in previous work.

We partition the 105 channels into ten groups based on their spatial location under the guidance of functional brain regions (Table 1). Specifically, we compared the spatial distribution of 105 electrodes with the spatial distribution of functional brain regions and mapped each electrode to its closest brain region. Details of the electrode spatial distribution can be found in (Hollenstein et al., 2018).

After the partition of the electrodes, we create a multi-view transformer model including ten convolutional transformers at the bottom level, where each convolutional transformer encodes the EEG signals from the electrodes in that region. On top of the ten convolutional transformers, we add a global transformer to unify the information from different brain regions. The combined information from the global transformer is further fed into the BART model for text decoding.

In summary, the multi-view transformer envisions multiple parallel convolutional transformer models where each captures different aspects of EEG signals combined from different spatial regions of the brain regions. This approach enhances the spatial resolution of the model and further improves the text decoding performance.

## 4. Experiment

### 4.1. Experimental Setup

**Dataset** We utilize both the Zuco (Hollenstein et al., 2018) and Image-EEG (Gifford et al., 2022) datasets for pre-training and use Zuco to train the multi-view transformer and BART model for text decoding. Details of both datasets are listed below.

- **Zuco** (Hollenstein et al., 2018) contains EEG and eye-tracking data from 12 healthy adult native English speakers engaged in natural English text reading for 4 - 6 hours. This dataset covers two standard reading tasks and a task-specific reading task, offering EEG and eye-tracking data for 21,629 words across 1,107 sentences and 154,173 fixations. Zuco contains both word-level EEG signals and sentence-level EEG signals. Sentence-level EEG refers to the complete original EEG signals recorded while the sub-

ject is reading some texts on the screen. The word-level EEG is generated based on sentence-level EEG, combined with eye-tracking data captured from the subjects. Specifically, the eye-tracking machine captures the coordinates of the screen where the subject’s gaze is focused, while also recording the current time. Then, the word corresponding to the captured coordinates on the screen is extracted.

- **Image-EEG** (Gifford et al., 2022) is a large and rich dataset containing high temporal resolution EEG signals of images of objects on natural backgrounds. The dataset included 10 participants, each performing 82,160 trials across 16,740 image conditions.

**Baselines** We compare EEG2TEXT with two baseline models for open-vocabulary EEG-to-text translation.

- **Baseline (EEGtoText)** (Wang & Ji, 2021) uses word-level EEG signals as input to a transformer model followed by a pre-trained BART model for decoding. EEGtoText is the first paper that proposed the open-vocabulary EEG-to-text translation task.
- **DeWave** (Duan et al., 2023) introduces a discrete codex encoding after the transformer layer, and uses both word-level EEG features and the raw EEG signals as input. DeWave is the most recent related work and it only included EEGtoText (Wang & Ji, 2021) as its baseline.

After a comprehensive literature review, we believe we have included all the baselines to our knowledge.

**Evaluation Metrics** We utilize BLEU-1, BLEU-2, BLEU-3, BLEU-4, and ROUGE-1 evaluation metrics to compare the performance of EEG2TEXT with the baselines.

The BLEU-N scores ( $N = 1, 2, 3, 4$ ) are used to measure the quality of the generated text, with higher values indicating better performance.

$$\text{BLEU} = \text{BP} \cdot \exp \left( \sum_{n=1}^N w_n \cdot \log \left( \frac{\text{count}_{\text{clip},n}}{\text{count}_{\text{ref},n}} \right) \right), \quad (3)$$

where BLEU represents the BLEU score; BP represents the brevity penalty;  $N$  represents the max n-gram order;  $w_n$  represents the n-gram weights;  $\text{count}_{\text{clip},n}$  represents count of candidate n-grams in reference and  $\text{count}_{\text{ref},n}$  represents count of reference n-grams.

ROUGE-1 scores, which include F (F1-score), P (precision), and R (recall), are used to evaluate the overlap between generated text and reference text.

$$\text{ROUGE-1} = \frac{\sum_{\text{ref}} \sum_{1\text{-gram}} \min(\text{match}, \text{ref})}{\sum_{\text{ref}} \sum_{1\text{-gram}} \text{ref}}, \quad (4)$$

Table 2. Optimal hyper-parameters for EEG2TEXT ablations.

Methods	Batch Size	Learning Rate
EEG2TEXT (Convolutional Transformer)	4	$1 \times 10^{-5}$
EEG2TEXT (+ Pre-training)	4	$5 \times 10^{-5}$
EEG2TEXT (+ Multi-View Transformer)	4	$3 \times 10^{-5}$

where ROUGE-1 represents the ROUGE-1 score; match represents the count of matching 1-gram; ref represents the count of 1-gram.

**Parameter Study** We used four A40 GPUs as our computing infrastructure and each training epoch took about 40 minutes. The optimal hyper-parameters for our results are listed in Table 2. The value ranges of each hyper-parameter are listed below:

- Batch Size  $\in \{4, 8, 16\}$
- Learning Rate  $\in \{1 \times 10^{-6}, 3 \times 10^{-6}, 5 \times 10^{-6}, 7.5 \times 10^{-6}, 8 \times 10^{-6}, 9 \times 10^{-6}, 1 \times 10^{-5}, 2 \times 10^{-5}, 3 \times 10^{-5}, 4 \times 10^{-5}, 5 \times 10^{-5}, 7.5 \times 10^{-5}, 1 \times 10^{-4}, 3 \times 10^{-4}, 5 \times 10^{-4}, 7.5 \times 10^{-4}, 1 \times 10^{-3}\}$
- Epoch  $\in \{15\}$

## 4.2. Results

**Main Results** Table 3 shows our main experimental results. The baseline method (Wang & Ji, 2021) achieves a moderate performance in text decoding with BLEU scores. DeWave (Duan et al., 2023) slightly improved the performance across all metrics, demonstrating the effectiveness of discrete encoding. EEG2TEXT improved the text decoding performance by a large margin due to several technical innovations. First, a single convolutional transformer achieved slightly lower BLEU scores (BLEU-1: -1.3%; BLEU-2: -0.5%; BLEU-3: -0.2%; BLEU-4: -0.0%) but higher ROUGE-1 scores (F1-score: +3.7%; Precision: +2.4%; Recall: -0.9%) compared to DeWave. Second, EEG2TEXT with pre-training further enhanced the BLEU scores (BLEU-1: +1.8%; BLEU-2: +1.9%; BLEU-3: +1.8%; BLEU-4: +1.6%) and ROUGE-1 scores (F1-score: +4.2%; Precision: +2.4%; Recall: +0.0%) compared to DeWave. Pre-training proved effective in enhancing text generation by providing a strong initialization foundation for our model. Third, EEG2TEXT with multi-view transformers achieved the highest scores across all metrics, with a significant increase in the BLEU scores (BLEU-1: +3.9%; BLEU-2: +5.0%; BLEU-3: +5.8%; BLEU-4: +5.9%) and ROUGE-1 scores (F1-score: +5.4%; Precision: +3.2%; Recall: +1.4%) compared to DeWave. EEG2TEXT excelled in generating coherent, contextually relevant, and high-quality text.

Table 3. Performance comparison of EEG2TEXT with baseline methods.

Methods	BLEU-N				ROUGE-1		
	N = 1	N = 2	N = 3	N = 4	F	P	R
Baseline (Wang & Ji, 2021)	0.401	0.231	0.125	0.068	0.301	0.317	0.288
DeWave (Duan et al., 2023)	0.413	0.241	0.139	0.082	0.288	0.337	0.306
EEG2TEXT (Convolutional Transformer)	0.400	0.236	0.137	0.082	0.325	0.361	0.297
EEG2TEXT (+ Pre-training)	0.445	0.274	0.175	0.117	0.341	0.383	0.310
EEG2TEXT (+ Multi-View Transformer)	<b>0.452</b>	<b>0.291</b>	<b>0.197</b>	<b>0.141</b>	<b>0.342</b>	<b>0.369</b>	<b>0.320</b>

Table 4. Ablation study of different input formats of the EEG signals.

Methods	BLEU-N				ROUGE-1		
	N = 1	N = 2	N = 3	N = 4	F	P	R
Spectrogram + Transformer	0.386	0.220	0.121	0.067	0.306	0.342	0.306
Spectrogram + Convolutional Transformer	0.374	0.209	0.112	0.061	0.302	0.339	0.274
EEG signal + Convolutional Transformer	<b>0.400</b>	<b>0.236</b>	<b>0.137</b>	<b>0.082</b>	<b>0.325</b>	<b>0.361</b>	<b>0.297</b>

**Convolutional Transformer** We first compare different input representations of the EEG signals to see how the representation affects the performance of a base convolutional transformer model. In this ablation study, we compare the raw EEG signals with their spectrograms using the fast Fourier transform (Cochran et al., 1967) to convert the original one-dimensional time array into a two-dimensional time-frequency matrix. The results are shown in Table 4. Using the raw EEG as the input consistently led to better performance than using the spectrogram as the input. Because the spectrogram only keeps the magnitude information and ignores the phase information of the raw EEG signal, the superior performance of the raw EEG signal suggested that the phase information might be important for decoding. Therefore, the raw EEG signals are used as the input in our subsequent experiments.

**EEG Pre-Training** We then conducted ablation experiments to compare the effectiveness of three pre-training strategies: 1) Masked Token Prediction (Devlin et al., 2018), 2) Continuous Masked Token Prediction, and 3) Re-Masked Token Prediction (Liu et al., 2019). The results are shown in Table 5. The Re-Masked Token Prediction (Liu et al., 2019) exhibits the best performance among all the three masking strategies. One potential reason is that the convolutional transformer model can learn more diverse semantic information by masking different tokens in each training epoch during pre-training.

In the above study, we focused on identifying the optimal pre-training strategy among the three without incorporating image-EEG data (Gifford et al., 2022). As an additional component, we introduced image-EEG data to assess the compatibility of our model with EEG signals from

multi-modal inputs. Leveraging our self-supervised pre-training strategy, we directly incorporated image-EEG data into the pre-training phase to enable the model to glean knowledge from diverse sources. The results, detailed in Table 6, demonstrate that adding image-EEG data significantly enhances translation performance for both the single convolutional transformer and the multi-view transformer.

**Multi-View Transformer** Finally, we compare different training strategies of the multi-view transformer to demonstrate the effectiveness of the multi-view transformer and find the best training strategy. The image-EEG data was not included in this ablation study. Specifically, we compared three training strategies as follows:

- **Only Global Transformer:** Fixing the parameters of all 10 convolutional transformer modules and training only the global transformer for text decoding.
- **Global Transformer + One Convolutional Transformer:** During each training epoch, randomly activate and train one convolutional transformer with the global transformer while fixing the parameters of the remaining nine convolutional transformers.
- **Global Transformer + Three Convolutional Transformers:** During each training epoch, randomly activate and train three convolutional transformers with the global transformer while fixing the parameters of the remaining seven convolutional transformers.

We have a large dataset with 2K batches to ensure each individual Transformer is trained sufficiently.

The results in Table 7 demonstrate that activating three convolutional transformers together with the global transformer

Table 5. Ablation study of different pre-training strategies of the EEG signals.

Methods	BLEU-N				ROUGE-1		
	N = 1	N = 2	N = 3	N = 4	F	P	R
Masked Token Prediction	0.409	0.242	0.141	0.087	0.325	0.357	0.300
Continuous Masked Token Prediction	0.411	0.243	0.137	0.078	0.319	0.352	0.294
Re-Masked Token Prediction	<b>0.431</b>	<b>0.260</b>	<b>0.157</b>	<b>0.098</b>	<b>0.330</b>	<b>0.361</b>	<b>0.306</b>

Table 6. Ablation study of adding image-EEG data into pre-training.

Methods	BLEU-N				ROUGE-1		
	N = 1	N = 2	N = 3	N = 4	F	P	R
Single-View without image-EEG	0.431	0.260	0.157	0.098	0.330	0.361	0.306
Single-View with image-EEG	0.445	0.274	0.175	0.117	0.341	0.383	0.310
Multi-View without image-EEG	0.442	0.277	0.179	0.121	0.335	0.365	0.311
Multi-View with image-EEG	<b>0.452</b>	<b>0.291</b>	<b>0.197</b>	<b>0.141</b>	<b>0.342</b>	<b>0.369</b>	<b>0.320</b>

achieves the best performance. This suggests further improvement may be attainable by increasing the number of activated convolutional transformers during each training epoch if more GPU resources are available.

**Case Study** Table 8 shows our case study results. In the first sentence, the baseline model accurately translates "good," whereas EEG2TEXT, in addition, accurately captures the first half of the sentence with "movie" (synonymous with "film"). Additionally, EEG2TEXT correctly translates the second half of the sentence with "disaster movie" corresponding to "monstrous one" in the original sentence. In the second sentence, EEG2TEXT accurately captured "won Nobel Prize in Chemistry," while the baseline produced incorrect information, stating "Pulitzer Prize" and the wrong field, "Literature." In the third sentence, both EEG2TEXT and the baseline correctly identified "book" and "Pulitzer Prize." However, EEG2TEXT, in addition, correctly identified the field as "Biography," while the baseline erroneously outputted "Fictionography."

In addition, we conducted an interesting case study to show that EEG2TEXT has the ability of zero-shot image-to-text translation. Details can be found in Appendix B.

## 5. Related Work

**Brain Computer Interface** The landscape of brain-to-speech and brain-to-text decoding encompasses three principal approaches grounded in the features they capture: motor imagery-based, overt speech-based, and inner speech-based. These methods explore a variety of brain signals, including electroencephalogram (EEG), electrocorticography (ECoG), and functional magnetic resonance imaging (fMRI). Despite these endeavors, existing approaches exhibit limitations con-

cerning vocabulary size, articulation dependence, speed, and device compatibility. Motor imagery-base systems, exemplified by point-and-click (Pandarinath et al., 2017) mechanisms and imaginary handwriting (Willett et al., 2021), show high accuracy but modest typing rates. Overt speech-based techniques for decoding speech offer expedited communication rates. However, they require either physical vocal tract movement (Herff et al., 2015; Anumanchipalli et al., 2019; Makin et al., 2020) or mental articulation imagination (Moses et al., 2021; Willett et al., 2023). This engenders language dependency and pronunciation variations across languages. Another line of research tackles articulation dependency by decoding imagined speech (Nieto et al., 2022) or reading text (Sun et al., 2019; Panachakel & Ramakrishnan, 2021). Our work follows this line of decoding reading text directly from EEG signals.

**EEG-to-Text Decoding** Prior investigations into the decoding of EEG-to-text, as documented in the literature (Herff et al., 2015; Sun et al., 2019; Anumanchipalli et al., 2019; Makin et al., 2020; Panachakel & Ramakrishnan, 2021; Moses et al., 2021; Nieto et al., 2022), have demonstrated commendable accuracy when applied to limited and closed vocabularies. Nevertheless, these studies encounter challenges in attaining comparable levels of accuracy when confronted with more extensive and open vocabularies. New investigations have expanded their focus from closed-vocabulary EEG-to-text decoding to encompass open-vocabulary scenarios (Wang & Ji, 2021; Willett et al., 2023; Tang et al., 2023; Duan et al., 2023). The two research studies most similar to our work are a baseline method (Wang & Ji, 2021) and DeWave (Duan et al., 2023). The baseline method proposes a framework utilizing transformer and pre-trained BART language models, which establish baseline performance of open-vocabulary EEG-to-

Table 7. Ablation study of different training strategies of the multi-view transformer.

Methods	BLEU-N				ROUGE-1		
	N = 1	N = 2	N = 3	N = 4	F	P	R
Only Global Transformer	0.404	0.238	0.139	0.084	0.303	0.335	0.279
+ One Convolutional Transformer	0.436	0.270	0.168	0.110	0.327	0.363	0.299
+ Three Convolutional Transformers	<b>0.442</b>	<b>0.277</b>	<b>0.179</b>	<b>0.121</b>	<b>0.335</b>	<b>0.365</b>	<b>0.311</b>

Table 8. Case study of the output sentences comparing EEG2TEXT and the baseline method (Wang &amp; Ji, 2021).

(1)	Ground Truth: It’s not a particularly <b>good film</b> , but neither is it a <b>monsterous</b> one.
	Baseline Output: was a a bad <b>good</b> story, but it is it <b>bad bad</b> . one.
	EEG2TEXT output: ’s a a <b>great</b> romantic <b>movie</b> , but it is it the <b>disaster</b> movie one.
(2)	Ground Truth: He won a <b>Nobel Prize in Chemistry</b> in 1928
	Baseline Output: was the Pulitzer Prize for Literature in 18.
	EEG2TEXT Output: won <b>Nobel Prize in Chemistry</b> for 1901
(3)	Ground Truth: The book was awarded the 1957 <b>Pulitzer Prize for Biography</b> .
	Baseline Output: first is published the Pulitzer <b>Pulitzer Prize</b> for Fictionography.
	EEG2TEXT Output: book is a <b>Pulitzer Prize for Biography</b> .

text translation. DeWave employs a quantization encoder to derive discrete encoding and aligns it with a pre-trained language model for the open-vocabulary EEG-to-text translation. The limitations of both the baseline method and DeWave lie in their reliance on eye-tracking calibration for word-level EEG feature extraction that introduces error propagation and lacks generalizability to scenarios such as inner speech decoding. EEG2TEXT improves the open-vocabulary EEG-to-text translation performance as well as enhancing the generality by requiring only sentence-level EEG signals as input.

**EEG Encoding** It is a challenging problem to effectively encode the long and noisy EEG signals to facilitate subsequent decoding tasks. In Conformer (Song et al., 2022), the authors propose a compact convolutional transformer, named EEG Conformer, to encapsulate local and global features in a unified EEG classification framework. Specifically, the convolution module learns the low-level local features throughout the one-dimensional temporal and spatial convolution layers. The self-attention module is straightforwardly connected to extract the global correlation within the local temporal features. However, in the case of the Conformer model, the authors trained this model from scratch, whereas EEG2TEXT further incorporated pre-training and multi-view settings to enhance the text translation performance.

**EEG Pre-Training** Recent work, such as BrainBERT (Wang et al., 2023), BENDR (Kostas et al., 2021) and

MAEEG (Chien et al., 2022), has been done on EEG signal pre-training that greatly inspired EEG2TEXT.

BrainBERT converts intracranial recordings to spectrograms and uses spectrograms as input. BrainBERT masks multiple continuous bands of random frequencies and time intervals from spectrograms and aims to reconstruct the original spectrogram. BENDR uses raw EEG signals as input. After a convolutional layer, the raw EEG signals are converted to embedding features. These embedding features are masked by using masked token prediction (Devlin et al., 2018) and the reconstruction goal is the original embedding features. MAEEG uses raw EEG signals as input and masks the embedding features of the convolutional layer generated with a masked token prediction as BENDR. However, MAEEG’s reconstruction goal is the raw EEG signals. EEG2TEXT directly masks the raw EEG signals with the pre-training objective to reconstruct the raw EEG signals. EEG2TEXT also experimented with various masking strategies and incorporated EEG signals for the pre-training process.

## 6. Conclusion

In this work, we proposed a novel EEG-to-text decoding model, EEG2TEXT that takes raw EEG signals as input and leverages EEG pre-training and a multi-view transformer to enhance the decoding performance. EEG2TEXT achieved superior performance for open-vocabulary EEG-to-text decoding. Future work includes expanding the model’s capabilities to EEG signals from diverse multi-modal data.



## 7. Acknowledgements

Our work is sponsored by the NSF NAIRR Pilot and PSC Neocortex, Commonwealth Cyber Initiative, Children’s National Hospital, Fralin Biomedical Research Institute (Virginia Tech), Sanghani Center for AI and Data Analytics (Virginia Tech), Virginia Tech Innovation Campus, and a generous gift from the Amazon + Virginia Tech Center for Efficient and Robust Machine Learning.

## References

- Aflalo, T., Kellis, S., Klaes, C., Lee, B., Shi, Y., Pejsa, K., Shanfield, K., Hayes-Jackson, S., Aisen, M., Heck, C., et al. Decoding motor imagery from the posterior parietal cortex of a tetraplegic human. *Science*, 348(6237):906–910, 2015.
- Anumanchipalli, G. K., Chartier, J., and Chang, E. F. Speech synthesis from neural decoding of spoken sentences. *Nature*, 568(7753):493–498, 2019.
- Bouton, C. E., Shaikhouni, A., Annetta, N. V., Bockbrader, M. A., Friedenberg, D. A., Nielson, D. M., Sharma, G., Sederberg, P. B., Glenn, B. C., Mysiw, W. J., et al. Restoring cortical control of functional movement in a human with quadriplegia. *Nature*, 533(7602):247–250, 2016.
- Chien, H.-Y. S., Goh, H., Sandino, C. M., and Cheng, J. Y. Maeeg: Masked auto-encoder for eeg representation learning, 2022.
- Cochran, W. T., Cooley, J. W., Favon, D. L., Helms, H. D., Kaenel, R. A., Lang, W. W., Maling, G. C., Nelson, D. E., Rader, C. M., and Welch, P. D. What is the fast fourier transform? *Proceedings of the IEEE*, 55(10):1664–1674, 1967.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL <http://arxiv.org/abs/1810.04805>.
- Duan, Y., Zhou, C., Wang, Z., Wang, Y.-K., and teng Lin, C. Dewave: Discrete encoding of EEG waves for EEG to text translation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=WaLI8slhLw>.
- Gifford, A. T., Dwivedi, K., Roig, G., and Cichy, R. M. A large and rich eeg dataset for modeling human visual object recognition. *NeuroImage*, 264:119754, 2022. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2022.119754>. URL <https://www.sciencedirect.com/science/article/pii/S1053811922008758>.
- Herff, C., Heger, D., De Pestere, A., Telaar, D., Brunner, P., Schalk, G., and Schultz, T. Brain-to-text: decoding spoken phrases from phone representations in the brain. *Frontiers in neuroscience*, 9:217, 2015.
- Hochberg, L. R., Bacher, D., Jarosiewicz, B., Masse, N. Y., Simeral, J. D., Vogel, J., Haddadin, S., Liu, J., Cash, S. S., Van Der Smagt, P., et al. Reach and grasp by people with tetraplegia using a neurally controlled robotic arm. *Nature*, 485(7398):372–375, 2012.
- Hollenstein, N., Rotsztein, J., Troendle, M., Pedroni, A., Zhang, C., and Langer, N. Zuco, a simultaneous eeg and eye-tracking resource for natural sentence reading. *Scientific data*, 5(1):1–13, 2018.
- Joshi, M., Chen, D., Liu, Y., Weld, D. S., Zettlemoyer, L., and Levy, O. Spanbert: Improving pre-training by representing and predicting spans. *CoRR*, abs/1907.10529, 2019. URL <http://arxiv.org/abs/1907.10529>.
- Kostas, D., Aroca-Ouellette, S., and Rudzicz, F. Bendr: using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers in Human Neuroscience*, 15:653659, 2021.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., and Zettlemoyer, L. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019. URL <http://arxiv.org/abs/1907.11692>.
- Lorach, H., Galvez, A., Spagnolo, V., Martel, F., Karakas, S., Intering, N., Vat, M., Faivre, O., Harte, C., Komi, S., et al. Walking naturally after spinal cord injury using a brain–spine interface. *Nature*, pp. 1–8, 2023.
- Makin, J. G., Moses, D. A., and Chang, E. F. Machine translation of cortical activity to text with an encoder–decoder framework. *Nature neuroscience*, 23(4):575–582, 2020.
- Martin, S., Iturrate, I., Millán, J. d. R., Knight, R. T., and Pasley, B. N. Decoding inner speech using electrocorticography: Progress and challenges toward a speech prosthesis. *Frontiers in neuroscience*, 12:422, 2018.
- Moses, D. A., Metzger, S. L., Liu, J. R., Anumanchipalli, G. K., Makin, J. G., Sun, P. F., Chartier, J., Dougherty, M. E., Liu, P. M., Abrams, G. M., et al. Neuroprosthesis for decoding speech in a paralyzed person with

- anarthria. *New England Journal of Medicine*, 385(3): 217–227, 2021.
- Nagel, S. and Spüler, M. Modelling the brain response to arbitrary visual stimulation patterns for a flexible high-speed brain-computer interface. *PloS one*, 13(10): e0206107, 2018.
- Nalborczyk, L., Grandchamp, R., Koster, E. H., Perrone-Bertolotti, M., and Lœvenbruck, H. Can we decode phonetic features in inner speech using surface electromyography? *PloS one*, 15(5):e0233282, 2020.
- Nieto, N., Peterson, V., Rufiner, H. L., Kamienkowski, J. E., and Spies, R. Thinking out loud, an open-access eeg-based bci dataset for inner speech recognition. *Scientific Data*, 9(1):52, 2022.
- Panachakel, J. T. and Ramakrishnan, A. G. Decoding covert speech from eeg—a comprehensive review. *Frontiers in Neuroscience*, 15:392, 2021.
- Pandarath, C., Nuyujukian, P., Blabe, C. H., Soric, B. L., Saab, J., Willett, F. R., Hochberg, L. R., Shenoy, K. V., and Henderson, J. M. High performance communication by people with paralysis using an intracortical brain-computer interface. *Elife*, 6:e18554, 2017.
- Song, Y., Zheng, Q., Liu, B., and Gao, X. Eeg conformer: Convolutional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:710–719, 2022.
- Sun, J., Wang, S., Zhang, J., and Zong, C. Towards sentence-level brain decoding with distributed representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 7047–7054, 2019.
- Tang, J., LeBel, A., Jain, S., and Huth, A. G. Semantic reconstruction of continuous language from non-invasive brain recordings. *Nature Neuroscience*, pp. 1–9, 2023.
- Wang, C., Subramaniam, V., Yaari, A. U., Kreiman, G., Katz, B., Cases, I., and Barbu, A. Brainbert: Self-supervised representation learning for intracranial recordings. *arXiv preprint arXiv:2302.14367*, 2023.
- Wang, Z. and Ji, H. Open vocabulary electroencephalography-to-text decoding and zero-shot sentiment classification. *CoRR*, abs/2112.02690, 2021. URL <https://arxiv.org/abs/2112.02690>.
- Willett, F. R., Avansino, D. T., Hochberg, L. R., Henderson, J. M., and Shenoy, K. V. High-performance brain-to-text communication via handwriting. *Nature*, 593(7858): 249–254, 2021.
- Willett, F. R., Kunz, E. M., Fan, C., Avansino, D. T., Wilson, G. H., Choi, E. Y., Kamdar, F., Glasser, M. F., Hochberg, L. R., Druckmann, S., et al. A high-performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.

## A. EEG to Spectrogram

Figure A1 shows a piece of EEG signals and its corresponding spectrogram.

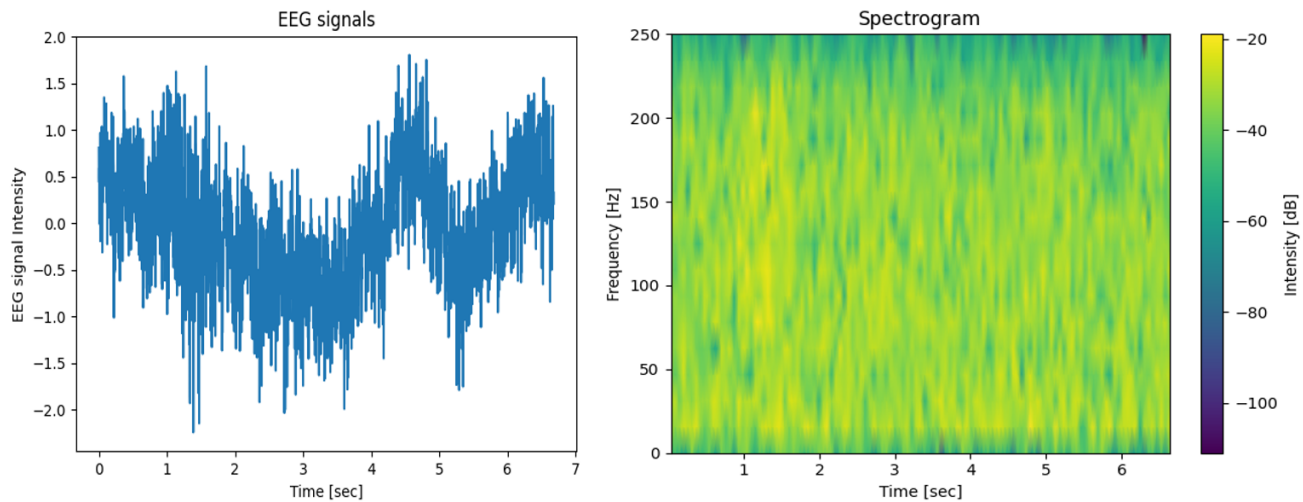


Figure A1. a piece of EEG signals and its corresponding Spectrogram

## B. Zero-Shot Image-to-Text Translation

Figure 2(a) and 2(b) show the zero-shot image-to-text translation results. We directly input the EEG signals of image-EEG data into the multi-view transformer model after training, and the output results are image-to-text translation results. The first image contains multiple cars, and the output accurately captures the "car" keyword. The second image contains a fish, and the output captures the "fish" keyword equally accurately.



(a) An image of car. The translation result of EEG2TEXT is: "alog,, **car**,,,,,,,,,,,,,,,,,,,,,,"



(b) An image of car. The translation result of EEG2TEXT is: "**fish**,,,,, has,,,,,,,,,,,,,,,,,,,,,,,,,,,,,"

Figure A2. Zero-Shot Image-to-Text Translation.