
MAD: Multi-Alignment MEG-to-Text Decoding

Yiqian Yang^{1*} Hyejeong Jo^{2*} Yiqun Duan^{3*} Qiang Zhang¹ Jinni Zhou¹
Won Hee Lee^{2†} Renjing Xu^{1†} Hui Xiong^{1†}

Abstract

Deciphering language from brain activity is a crucial task in brain-computer interface (BCI) research. Non-invasive cerebral signaling techniques including electroencephalography (EEG) and magnetoencephalography (MEG) are becoming increasingly popular due to their safety and practicality, avoiding invasive electrode implantation. However, current works under-investigated three points: 1) a predominant focus on EEG with limited exploration of MEG, which provides superior signal quality; 2) poor performance on unseen text, indicating the need for models that can better generalize to diverse linguistic contexts; 3) insufficient integration of information from other modalities, which could potentially constrain our capacity to comprehensively understand the intricate dynamics of brain activity. This study presents a novel approach for translating MEG signals into text using a speech-decoding framework with multiple alignments. Our method is the first to introduce an end-to-end multi-alignment framework for totally unseen text generation directly from MEG signals. We achieve an impressive BLEU-1 score on the *GWilliams* dataset, significantly outperforming the baseline from 5.49 to 10.44 on the BLEU-1 metric. This improvement demonstrates the advancement of our model towards real-world applications and underscores its potential in advancing BCI research. Code is available at <https://github.com/NeuSpeech/MAD-MEG2text>.

1 Introduction

Decoding brain to language has emerged as a rapidly developing area of neurotechnology, offering semantic communication and control for general Brain-Computer-Interface (BCI) tasks. This region has garnered growing focus as it may profoundly impact individuals with verbal and movement disabilities resulting from conditions such as severe spinal cord trauma or end-stage amyotrophic lateral sclerosis (ALS). Moreover, the scope of brain-to-text technology extends to pioneer novel human-machine interfaces, allowing seamless control of prosthetic limbs, software, and virtual environments, shifting the paradigm of interaction for both able-bodied individuals and those with disabilities, and re-defining what is achievable in both everyday life and professional spheres.

*These authors contributed equally to this work

†These are corresponding authors

¹The Hong Kong University of Science and Technology (Guangzhou), People's Republic of China,

Yiqian Yang and Hui Xiong are with AI Thrust, HKUST(GZ).

Qiang Zhang and Renjing Xu are with MICS Thrust, HKUST(GZ),

Jinni Zhou is with RBM Base, College of Future Technology,

Email: yyang937@connect.hkust-gz.edu.cn, xionghui@hkust-gz.edu.cn, qzhang749@connect.hkust-gz.edu.cn,

renjingxu@hkust-gz.edu.cn, eejinni@hkust-gz.edu.cn

²Department of Software Convergence, Kyung Hee University, Republic of Korea,

Email:girlsending0@khu.ac.kr, whlee@khu.ac.kr

³GrapheneX-UTS HAI Centre, Australia Artificial Intelligence Institute, University of Technology Sydney, Australia

Email:duanyiquncc@gmail.com

Under this scope, various previous works have explored this area in multiple ways. Pioneer researchers first verify this idea by using invasive signals such as Electrocorticography (ECoG) [1, 2, 3, 4]. Recently, these invasive methods [5, 6] concentrate on decoding speech, phonemes or letter from ECoG signals and have achieved remarkably high accuracy using limited word sets for real-time brain-to-text translation. However, these invasive-signal-based approaches pose significant medical risks and challenges for long-term use.

Non-invasive techniques, therefore, present a safer and more sustainable alternative, albeit with their own set of challenges. Wang et al. [7] showcased a method for translating EEG signals into text with an extensive lexicon, utilizing language models that had been pre-trained on EEG data features at word-level. Duan et al. [8] progressed this methodology by interpreting raw EEG signals directly, devoid of reliance on temporal indicators, but their models still relied heavily on teacher forcing for evaluation, limiting their ability to generate meaningful sentences autonomously in real-life scenarios. At the same time, although Magnetoencephalography (MEG) provides better signal quality, previous works [9, 10, 11] on MEG have primarily focused on decoding limited classes or short phrases from MEG signals, showing limited success in generating whole sentences and complete semantic segments.

Furthermore, as pointed out by Jo et al. [12], all previous works in EEG-to-Text translation following Wang’s method [7] meets the “decoder dominated” problem. It means that given a strong decoder and noisy EEG input, these models are more likely to memorize the text distribution corresponding to certain statistical features rather than mapping EEG to semantic texts. Thus, these models have similar performances even when we replace EEG input with random noise. Besides, due to the nature of limited data and the non-understandability of the neural signal, it is difficult to train and evaluate the model. Yang et al. [13] proposed NeuSpeech model on MEG to text task, however, their model is evaluated on the text that is seen in the training set, which does not meet the need for open-vocabulary translation. Defossez et al. [14] highlighted the potential to decode speech perception from MEG signals, where they matched MEG signals with corresponding speech segments. However, their approach was limited to classification tasks and could not generate sentences directly from MEG signals. This underscores a significant gap in the current state of MEG-based brain-to-text decoding.

In this paper, our motivation is to establish an end-to-end framework for open-vocabulary MEG-to-Text translation capable of processing unseen text without relying on biomarkers, while ensuring that the encoder captures brain dynamics effectively. We propose Multi-Alignment MEG-to-Text Decoding (MAD) with the aim of guiding the brain encoders towards learning salient representations. To achieve this, we incorporate audio as an auxiliary modality to facilitate alignment. Here, we make a bold assumption that directly formatting noise brain signals into discrete text is difficult due to limited data. Hence, we utilize brain module [14] and an extra whisper model [15] to align brain representation in three aspects as shown in Figure 1, the Mel spectrogram, hidden state, and text. 1) We first align the Brain module with audio in the Mel spectrogram feature space to learn low-level features, such as acoustic features. 2) Additionally, we align the hidden state output from both the whisper encoder and the brain module in latent space, enhancing the model’s ability to extract high-level semantic features. 3) Lastly, we align the text representation from both streams within the framework.

Our objective in incorporating textual data is to assess whether it can furnish supplementary contextual cues that enhance the correspondence between neural activity and the resulting linguistic output.

Comprehensive experiments are conducted by utilizing non-invasive public MEG data from *GWilliams* [16] dataset, which captured MEG signals during a speech listening task. Remarkably, **MAD is capable of generalizing to unseen text**. Performance is evaluated using translation text relevancy metrics [17, 18]. On raw MEG waves, MAD achieves 10.44 BLEU-1 on *GWilliams* **without teacher-forcing** evaluation on **entirely unseen text** which largely exceeds the current SOTA performance. This paper also provides insights through numerous ablation studies to help people understand the impact of each component on aligning the MEG signal with texts. The contributions of this research could be summarized as follows:

- MAD presents an end-to-end neural network design for the direct conversion of MEG signals into text in open-vocabulary, obviating the dependence on markers, teacher forcing, or pre-training, representing the initial implementation of translating raw MEG waves into text for unseen content.

- We are the first to investigate various alignments and demonstrate the benefits of aligning with speech modality rather than text modality in the MEG-to-text transcription task, offering significant insights for network improvement.
- Our extensive experimentation and thorough analysis of the proposed model showcase its effectiveness and highlight its superiority over existing methods in terms of translation accuracy, efficiency, and reliability.

2 Related Works

The discipline of converting brain signals into textual output has undergone considerable development in the contemporary era. In 2019, Anumanchipalli et al. [1] introduced a pioneering model capable of translating ECoG patterns into the articulatory movements necessary for speech production, subsequently generating acoustic properties such as MFCCs, leading to the production of intelligible speech. This landmark study ignited further exploration within the field. In the subsequent year, Wang et al. [2] leveraged the capabilities of generative adversarial networks (GANs) to decipher ECoG data and synthesize speech. The year following, Willett et al. [3] engineered a system that utilized a recurrent neural network (RNN) alongside a probabilistic language model to decode letters from neural activity during the act of handwriting. Most recently, Metzger et al. [19] constructed a sequence of processes that converted ECoG signals into textual information using an RNN, enhancing the results with the GPT-2 language model.

Within the domain of open-vocabulary interpretation, Metzger et al. [6] unveiled an RNN architecture capable of real-time decoding of speech, text, sentiment, and facial expressions from ECoG data. Simultaneously, Willett et al. [5] managed to interpret text directly from neural activity. Liu et al. [20] introduced a tripartite model designed to decode logo-syllabic languages, such as Chinese, by transforming ECoG signals into Chinese pinyin inclusive of tones and syllables, followed by speech synthesis. In a related development, Feng et al. [21] achieved text interpretation from SEEG recordings. It is essential to highlight that these functional systems are predominantly reliant on invasive neural recordings.

In the domain of non-invasive neural recording, Meta unveiled a brain-to-speech system that leverages contrastive learning with MEG and EEG data [14]. While this system is proficient in categorizing a constrained set of sentences, it is not conducive to open-vocabulary textual interpretation. Ghazaryan et al. [11] explored the decoding of a restricted vocabulary from MEG responses. Wang et al. [7] crafted a mechanism for translating EEG features at the word level into text, employing a pre-trained BART model [22]. Subsequent investigations, including Dewave [8], adopted the methodology established by Wang et al. [7], proposing a schema that incorporates wave2vec [23] and discrete codex for robust representations, which are subsequently funneled into a BART [22] model for text synthesis. These approaches, however, are dependent on teacher-forcing and disregard the necessity of comparing results with noise-injected inputs, potentially resulting in an inflated assessment of system efficacy. Recent scholarship [12] has revealed the limitations of these methods.

Yang et al. [13] proposed an end-to-end paradigm for converting MEG signals to text, demonstrating high performance when training and evaluation sets were fully overlapped. However, it does not show good performance on unseen text. Our approach diverges from these methods by employing transfer learning with assistance of extra modality (Mel spectrogram) to align the model through multiple stages with low-level and high-level features of the ground truth. This enables our model to learn more effectively and generalize better to unseen text.

3 Method

3.1 Task Definition

Given a sequence of raw segment-level MEG signals ε , the goal is to decode the associated open-vocabulary text tokens T . This task also incorporates additional information in the form of speech Ξ . The MEG-Speech-Text pairs $\langle \varepsilon, \Xi, T \rangle$ are collected during speech perception. Our approach focuses on decoding T using only the raw MEG signal ε , with the support of Ξ . MAD represents the first attempt at tackling this MEG to unseen text translation challenge.

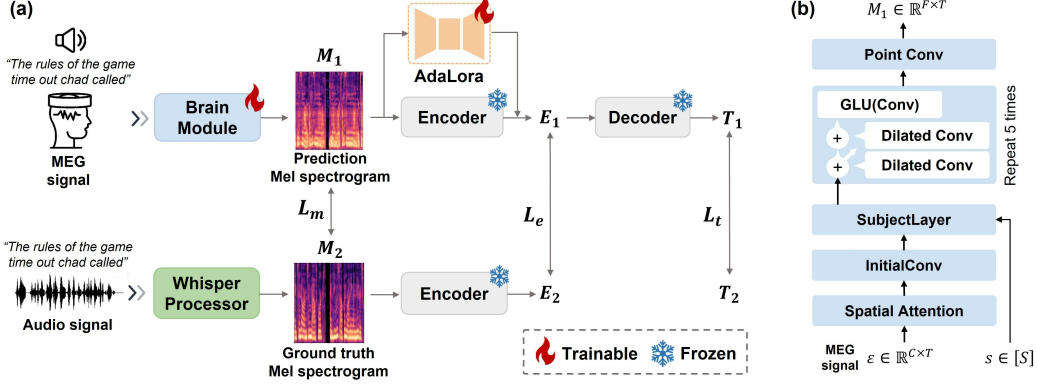


Figure 1. (a) Overview of model architecture. We added alignments on the Mel spectrogram, the hidden states, and the text. There are three types of alignment, which are either based on our physics world (text and speech) or a largely pre-trained model. M_1, M_2 is predicted and ground truth Mel spectrogram, E_1, E_2 is the hidden state of meg-input and speech-input encoder respectively. T_1 , and T_2 are predicted and ground truth text respectively. (b) Details about the brain module.

3.2 Model

Figure. 1 shows the overview of our work. Our model uses some transfer learning techniques to facilitate better performance on unseen text. The encoder and decoder models are from the Whisper model [15], a transformer-based encoder-decoder architecture tailored for robust speech recognition in challenging environments such as noisy conditions. The brain module [14] first takes the MEG signal ε of C channels in the Spatial Attention layer, it adds position embedding of physical sensors to the MEG, then Initial Conv maps the MEG channel number to hidden model dimensions. After that, the Subject Layer takes the MEG feature and subject index and applies subject embedding on the MEG feature. Next, the MEG feature is input into the residual-designed module which is repeated 5 times. Finally, after Point Conv of which kernel size is 1, it maps to the Mel spectrogram M_1 .

L_m is the loss to align the Mel spectrogram, which is Clip loss [24] in this situation. Then we want to make sure the encoder model can learn high-level features, so we designed to align the encoder output with L_e , which is Maximum Mean Discrepancy (MMD) loss [25]. We used LoRA [26] module to train our architecture for saving memory. Last but not least, we have the cross entropy loss L_t for predicted text and ground truth text. The overall loss L is below:

$$L = \lambda_m \cdot L_m + \lambda_e \cdot L_e + \lambda_t \cdot L_t \quad (1)$$

Recall the clip loss [24] function, it takes two feature representations from each modality. These features are then used to calculate the similarity scores between the representations of the image and text modalities. The Clip loss function aims to minimize the distance between matching pairs of image and text representations while maximizing the distance between non-matching pairs. This approach allows the CLIP model to learn a joint embedding space where semantically similar image-text pairs are close together, enabling tasks like zero-shot image classification and text-based image retrieval. Here the clip loss is applied on the Mel spectrogram, which is of 3 dimensions, so we flattened the batch size and time length dimensions as the first dimension, and then the loss is calculated as follows:

The MMD loss (Maximum Mean Discrepancy loss) is a measure of the discrepancy between two probability distributions. It is commonly used in domain adaptation and generative modeling to encourage the distributions of source and target data to be similar. If we flatten the hidden state E of the batch size n , time dimension t_d and feature dimension d_e , it will run out of memory if we input full length into the model, so we randomly select features time wise t_r , therefore the selected features is E_r shape is $[n, t_r, d_e]$ The formula for the MMD loss is:

$$\text{MMD}^2(E_1, E_2) = \frac{1}{n} \left\| \sum_{i=1}^n \phi(E_{1r}(i)) - \sum_{i=1}^n \phi(E_{2r}(i)) \right\|_{\mathcal{H}}^2 \quad (2)$$

Algorithm 1: CLIP-like Loss Calculation

Input: $M_1 [n, d_m]$ Predicted Mel spectrogram ,
 $M_2 [n, d_m]$ Ground truth Mel spectrogram ,
 d_m Dimensionality of multimodal embedding,
 t Learned temperature parameter,
 n Batch size.
Output: CLIP loss

```
1 logits  $\leftarrow M_1 \cdot M_2^T \cdot e^t$ ; // Scaled pairwise cosine similarities, [n,n]
2 labels  $\leftarrow \text{Range}(n)$ ; // Labels for each example
3 loss1  $\leftarrow \text{CrossEntropyLoss}(\text{logits}, \text{labels}, \text{axis} = 0)$ ;
4 loss2  $\leftarrow \text{CrossEntropyLoss}(\text{logits}, \text{labels}, \text{axis} = 1)$ ;
5  $L_m \leftarrow \text{Mean}(\text{loss}_1, \text{loss}_2)$ ;
6 return  $L_m$ ;
```

For an Automatic Speech Recognition (ASR) system, the cross-entropy loss is commonly used as a loss function to train the model. The basic idea is similar to the general cross-entropy loss but adapted for the ASR context where the inputs are speech features and the outputs are text transcriptions. The cross-entropy loss in the context of ASR can be defined as follows:

$$\text{CrossEntropyLoss} = -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \sum_{c=1}^C T_{1,i,t,c} \log(T_{2,i,t,c}) \quad (3)$$

3.3 Evaluation

We evaluate transcribing performance on *GWilliams* dataset [16] using NLP metrics, BLEU [17] is used to evaluate the accuracy of machine-translated text, ROUGE-1-F [18] is to measure the quality of automatic summarization, BertScore [27] is a measurement of semantic similarity, CER [28] is used to evaluate the accuracy of speech recognition and self-BLEU [29] is used to assess the diversity of generated text.

1. BLEU is used to evaluate the accuracy of machine-translated text.
2. ROUGE-1-F is to measure the quality of automatic summarization.
3. BertScore is a measurement of semantic similarity.
4. CER is used to evaluate the accuracy of speech recognition.
5. Self-BLEU is used to assess the diversity of generated text.

4 Experiments

4.1 Dataset

The GWilliams dataset [16] is a magnetoencephalography (MEG) dataset designed for assessing natural speech comprehension. It features authentic MEG recordings from 27 participants proficient in English. These participants engaged in two separate sessions, each involving two hours of listening to four stories, which are “cable spool fort”, “easy money”, “lw1”, “the black willow”. To get a fair evaluation, we split our dataset directly on stories, we test on “cable spool fort”, validate on “lw1” and train on other stories. Details are in Table. 1. For more details about the dataset, please refer to Supp. A.

For preprocessing, we used first band pass filter the MEG signal ε between 1 Hz and 40 Hz, then it is resampled to 100Hz to reduce computing. We ensure that we separated training, evaluation, testing set totally since we used one story for testing, another story for evaluation, last two ones for training. We extract 4-second windows from the MEG-speech-text pairs, sliding every second and randomly shifting the window by ± 0.5 seconds to generate samples. Speech Ξ is then transformed to Mel M with window length of 400, hop length of 160, which is the original configuration in Whisper

Table 1. Details about the dataset splits, we ensured the three splits are totally separated. Unique sentences means the sentences that are different with other sentences, same meaning for unique words. There is no overlap sentence between train and test set. 371(46%) means 371 words in test set is also in train set, accounting for 46 percentage.

Split	Segments	Unique sentences	Words	Unique words	Overlap sentence	Overlap words
train	133966	13266	150497	2776	-	-
validation	14896	1387	156027	478	-	-
test	31115	3151	355654	805	0	371(46%)

Table 2. Comparison with other models. Lo is LoRA, B is brain module. Bert here means Bertscore. Results is obtained without teacher forcing in evaluation. Here, Tr stands for trainable modules. B-1 stands for BLEU-1. R-1 stands for ROUGE-1-F. SB stands for Self-BLEU. RS means randomly selecting sentences from test set as predictions. As we can see, only MAD is much higher than RS on BLEU-1 score.

Modality	Method	Tr	Loss	B-1(%) \uparrow	R-1(%) \uparrow	Bert(%) \uparrow	CER(%) \downarrow	SB(%) \downarrow
-	RS	-	-	5.86	7.20	83.73	87.30	96.12
MEG	NeuSpeech [13]	Lo	L_t	5.49	8.43	83.98	77.02	99.7
MEG	Wav2vec2CTC [14]	B	L_m	0.55	1.44	76.02	152.23	92.67
MEG	MAD	B	$L_m + L_e$	10.44	6.93	83.39	89.82	85.66
Noise	MAD	B	$L_m + L_e$	3.87	3.16	83.20	126.95	87.54
MEG	MAD w/tf	B	$L_m + L_e$	12.93	18.28	82.87	74.31	83.35
Noise	MAD w/tf	B	$L_m + L_e$	0.19	6.68	59.92	87.57	68.63

model [15], since the setted speech sampling rate is 16kHz, after conversion, M is of shape [400, 80] time and feature wise for 4 second speech, then it is matched with ε of time length 400.

4.2 Implementation details

All models were trained using Nvidia 4090 (24GB) GPUs. Training was conducted with a learning rate of $3e-4$ and a batch size of 32 over 5 epochs, selecting the best-performing model based on evaluation loss. AdamW was employed as the optimizer across all models. Each experiment takes about 18 hours on signal GPU with 8 workers to finish. Lambda value in all experiment on MAD model set as follows: $\lambda_m = 1$, $\lambda_e = 0.01$, $\lambda_t = 1$.

4.3 Evaluation Metrics

The performance comparison of our proposed MAD model with other state-of-the-art models is summarized in Table 2. The table highlights various configurations and the corresponding evaluation metrics: BLEU-1, ROUGE-1, BertScore, CER and self-BLEU. Each model’s performance is evaluated on MEG data, with results illustrating the impact of different loss functions and modules on decoding accuracy.

We compare the performance of our proposed model, MAD, against existing state-of-the-art methods, NeuSpeech [13] and Wav2vec2CTC [14], for decoding MEG signals into text. The performance metrics used for evaluation include BLEU-1, ROUGE-1-F, Bertscore, and Character Error Rate (CER). The results are summarized in Table 2. We find out BLEU-1 seems to be the most effective measurement in this situation.

NeuSpeech [13] is a encoder-decoder framework model used for MEG, utilizing the Low-Rank Adaptation (LoRA) method with a text-based loss (L_t), achieves best scores on ROUGE-1-F, BertScore, and CER. However, the self-bleu score is almost 100%, which means the generation always repeat same thing. Besides, the BLEU-1 score is lower than RS, which means these three metrics are not reliable, which is further discussed in Supp. B.

Wav2vec2CTC [14]: The original model predicts the output of the Wav2vec2 [23] encoder with brain module. We add the pretrained language model head in the Wav2vec2CTC [23] model as another baseline. This model shows significantly lower performance across all metrics, which is not effective.

Our MAD model, which integrates the brain module with a combined loss ($L_m + L_e$), demonstrates superior performance with a BLEU-1 score of 10.44% which is about 5 points higher than NeuSpeech [13] and RS. Besides, we compared the performance of our model when it receives pure Gaussian noise which is the shape of the MEG signal to show that our model is generating text based on MEG signal. For noise input, MAD’s performance BLEU-1 dropped to 3.87%, indicating that MAD model has learned from the MEG signal rather than just noise. Additionally, we evaluated MAD with teacher forcing. When teacher forcing was applied (MAD w/tf), the model’s performance significantly improved, achieving a BLEU-1 score of 12.93% and a ROUGE-1-F score of 18.28%, confirming the effectiveness of teacher forcing in enhancing model performance. Similarly, the BLEU-1 score for noise w/tf is low too (0.19%), further indicating our model can distinguish noise and MEG. In addition, our model has low Self-BLEU which means our model is generate diverse sentences according to MEG signal rather than simply repeating.

Overall, our MAD model achieved state-of-the-art (SOTA) performance for MEG-to-text decoding compared to previous SOTA models, demonstrating significant progress in MEG-to-text translation. Additionally, we performed a fair comparison with noise and RS, which served as two error bars to validate the robustness and reliability of our model’s performance. Furthermore, the self-BLEU scores indicated the diversity of our model’s generated text, demonstrating its ability to truly learn and generalize from the data. Next section, we will show the generated sample along with the Mel spectrogram to further show the effectiveness of our MAD model.

4.4 Generated Samples

4.4.1 Text

Table 3. Transcription results. These are some results obtained with teacher forcing evaluation. **Bold** for exact matched words, underline for similar semantic words.

Decoding Results on <i>GWilliams</i> [16]
Ground Truth: corner of his eyes two forts stood on the playground and a hot
Prediction: the own of his <u>in the ground for</u> and the few spot hot
Ground Truth: of the top of the black fort like a <u>gold</u> headed monster tucker
Prediction: of the top hole a giant <u>medal</u> to is
Ground Truth: he knew it <u>wasn’t</u> going to be as <u>relax easy</u> as just pretending he was too
Prediction: he <u>wast</u> a to be a bad as as as it a to was a <u>lazy</u>

We showed the text result in table 3. It presents the transcription results obtained using the teacher forcing evaluation method. The transcription results indicate that while the model can generate segments of text that partially match the ground truth, there are significant gaps in overall accuracy and coherence. Specifically, for the first example, the ground truth is “corner of his eyes two forts stood on the playground and a hot” and the model’s prediction is “the own of his in the ground for and the few spot hot”. Although the model captures some keywords like “ground” and “hot,” the overall sentence diverges significantly from the ground truth, exhibiting repetition and grammatical errors. This outcome highlights the model’s struggle with complex sentence structures and semantic relationships.

In the second example, the ground truth is “of the top of the black fort like a gold headed monster tucker,” while the model predicts “of the top hole a giant medal to is.” Here, the model successfully identifies “of the top,” and the subsequent key word matches the ground truth in semantics, particularly “medal” which is semantically similar to “gold” in the ground truth sentence. Though it ignores the “monster tucker”, this suggests that, despite the MAD model’s failure in maintaining coherence and context understanding between words, it can yield some keywords which are semantically similar to the keywords in the original sentence.

For the third example, the ground truth is “he knew it wasn’t going to be as relax easy as just pretending he was too,” whereas the model predicts “he wast a to be a bad as as as it a to was a lazy.” Although the initial word “he” matches the ground truth, the following prediction includes repeated words and grammatical errors, making it difficult to form a meaningful sentence. However, we should notice that “lazy” may be similar to “relax” in meaning.

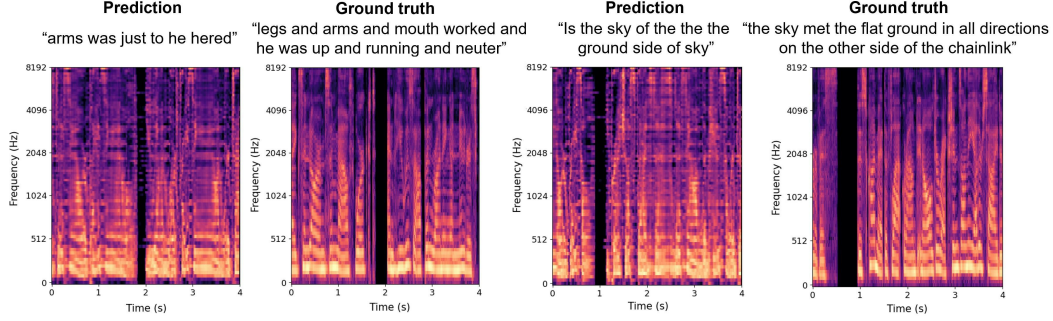


Figure 2. Two sample examples from the test set. Predictions refer to Mel spectrograms generated by the brain module. Ground truth refers to Mel spectrograms of the audio signal processed by the whisper processor. The predicted text was generated using teacher forcing.

For these three texts, we can observe that MAD can generate semantically similar words though they may not be coherent. Besides, the noise inputs generated only blanks. This demonstrates that our MAD model can capture the semantic meaning in the MEG signals, rather than only relying on the decoder.

Overall, these results reflect the inherent difficulties in directly decoding MEG signals into natural language text. While the model demonstrates some capability in recognizing individual words, there is substantial room for improvement in generating coherent and accurate sentences. Particularly, the model struggles with complex grammatical structures and longer sentences. These findings underscore the necessity of further optimizing the decoding model, especially in enhancing contextual awareness and semantic coherence. Future work should focus on improving the model’s ability to understand context and the relationships between words to enhance the overall accuracy and readability of the transcriptions.

4.4.2 Mel spectrogram

More than text, we showed the Mel spectrogram in Figure 2. It presents the Mel spectrogram of the two sample sentences in the test set. In this context, it is employed to compare the predicted audio signal generated by the model with the actual ground truth audio signal in the form of Mel spectrogram.

Upon examining the spectrograms of two samples, several observations can be made regarding the model’s capabilities and performance. 1) There is a general similarity between prediction and ground truth in the overall structure, 2) the model learns some fine-grained details such as temporal variations in the low-frequency regions which have bigger energy than the high-frequency region, 3) the model can predict the speech signal’s temporal blanks, proving it understands the MEG features associated with the absence of speech. However, significant discrepancies are apparent. While the ground truth spectrogram displays a more complex and detailed pattern with distinct frequency bands and variations over time, the predicted spectrogram seems less detailed and exhibits more uniform and repetitive patterns.

These discrepancies highlight the current limitations of the model in producing high-quality, accurate, natural audio signals from MEG data. Future work can introduce pre-trained generative models in speech modality to improve the model’s ability to learn and represent these fine-grained details, which is important for accurate speech recognition.

4.5 Model Ablation

Table 4 presents a comparison of various configurations of trainable modules and loss functions in the brain-to-text decoding model, evaluated under teacher forcing conditions. The configurations include different combinations of the brain module (B), LoRA applied to the encoder (Lo), the encoder (E), and the decoder (D). The evaluation metrics used are BLEU-1, ROUGE-1, Bert score, and CER (Character Error Rate), Self-BLEU.

Table 4. Here shows the comparison of using different modules and loss. B means brain module, Lo means LoRA applied on encoder. E means encoder, D means decoder. These results are obtained without teacher forcing in evaluation. $L_m(mmd)$ is the mmd loss for aligning mel spectrogram instead of Clip loss. B-1 is the abbreviation of BLEU-1. R-1 is the ROUGE-1-F. SB is self-BLEU.

Loss	Trainables	Architect.	B-1 (%)↑	R-1 (%)↑	Bert (%)↑	CER (%)↓	SB (%)↓
L_e	B	B+D	10.09	6.29	82.74	88.84	83.62
$L_e + L_t$	B	B+D	6.15	4.81	84.43	80.33	95.32
L_m	B	B+E+D	1.88	2.24	79.83	83.65	99.03
$L_m + L_e$	B	B+E+D	10.44	6.93	83.39	89.82	85.28
$L_m(mmd) + L_e$	B	B+E+D	9.64	5.71	81.62	87.95	80.55
$L_m + L_e + L_t$	B	B+E+D	7.14	4.37	82.29	88.40	83.95
$L_m + L_e$	B+Lo	B+E+D	1.13	0.79	81.17	87.65	99.98
$L_m + L_e + L_t$	B+Lo	B+E+D	8.33	6.40	83.14	91.43	99.11

The baseline configuration, using only the brain module with the loss function L_m , achieves a BLEU-1 score of 1.88 and a ROUGE-1 score of 2.24, with Bert and CER scores of 79.83 and 83.65, respectively. Self-BLEU scores of 99.03 indicate the model generated almost identical sentences, showing that using only the brain module results in significant errors and inaccurate content.

Adding the encoder loss L_e to L_m while maintaining the same modules (B+E+D) significantly improves performance, yielding a BLEU-1 score of 10.44 and a ROUGE-1 score of 6.93. The Bert and CER scores were 83.39 and 89.82, respectively. In this configuration, we changed the alignment loss from clip to MMD($L_m(mmd) + L_e$), resulting in BLEU-1 scores of 9.64 and ROUGE-1 scores of 5.71. Similarly, the configuration using the brain module with L_e (B+D) achieves the following scores: a BLEU-1 score of 10.09, a ROUGE-1 score of 6.29, a BERT score of 82.74, and a CER of 88.84. This indicates enhanced decoding accuracy when using encoder loss. Adding the triplet loss L_t to this configuration decreases the BLEU-1 and ROUGE-1 scores to 6.15 and 4.81, respectively.

Using LoRA with the combination of L_m and L_e results in significantly poor performance, with BLEU-1 and ROUGE-1 scores of 1.13 and 0.79, and Bert scores of 81.17. The Self-BLEU score of 99.98 indicates that this configuration is highly ineffective, likely due to an incompatibility between LoRA and the task requirements. Incorporating the triplet loss L_t along with L_m and L_e for the same architecture (B+E+D) resulted in a Self-BLEU score of 99.11, indicating that the model generated almost identical sentences.

The results indicate that the brain module (B) is crucial for effective brain-to-text decoding, and the combination of multiple loss functions, particularly with the inclusion of the encoder loss (L_e), enhances the model’s performance. Configurations involving LoRA applied to the encoder are generally less effective unless complemented with the L_t , highlighting the need for carefully designed adaptation strategies for optimal performance in this context.

5 Limitation

Although our MAD model outperforms previous SOTA models, we have to point out that this model’s generation is far from practical utilization in reality since the performance is much lower than speech recognition models. Besides, this work is implemented on listening datasets, which is different from silent speech.

6 Conclusion

In this paper, we presented MAD, a novel end-to-end training framework for MEG-to-Text translation. Our model leverages a multi-stage alignment utilizing auxiliary audio modality, which aligns brain activity data more effectively with corresponding textual outputs. Experimental results suggest that the newly proposed MAD framework achieves 10.44 BLEU-1 on *GWilliams* **without teacher-forcing** evaluation on **entirely unseen text** which largely exceeds the current SOTA performance. Through comprehensive ablation studies, we demonstrated the performance of our approach in various situations. Our results indicate that the brain module, in conjunction with appropriate loss functions, substantially enhances decoding performance. The inclusion of encoder and decoder

modules further refines the text generation process, with the triplet loss playing a crucial role in improving the model’s robustness and accuracy. Particularly, the combination of the brain module with both the encoder and decoder, enhanced by multiple loss functions, shows marked improvements in BLEU-1 and ROUGE-1 scores, while reducing word and character error rates. The insights gained from this research underline the potential of the MAD framework in the realm of neural decoding. By effectively capturing the complex patterns in MEG signals and translating them into coherent text, our approach offers a promising solution for brain-to-text applications. This work sets the stage for further exploration into multi-modal alignments and their impact on neural decoding systems.

In conclusion, our proposed MAD framework significantly advances the state-of-the-art in brain-to-text decoding, offering new avenues for enhancing communication tools for individuals with severe speech and motor impairments. Future work will focus on refining alignment mechanisms and extending the application of our model to more diverse linguistic tasks.

References

- [1] Gopala K. Anumanchipalli, Josh Chartier, and Edward F. Chang. Speech synthesis from neural decoding of spoken sentences. *Nature*, 568(7753):493–498, April 2019.
- [2] Ran Wang, Xupeng Chen, Amirhossein Khalilian-Gourtani, Zhaoxi Chen, Leyao Yu, Adeen Flinker, and Yao Wang. Stimulus speech decoding from human cortex with generative adversarial network transfer learning. In *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, April 2020.
- [3] Francis R. Willett, Donald T. Avansino, Leigh R. Hochberg, Jaimie M. Henderson, and Krishna V. Shenoy. High-performance brain-to-text communication via handwriting. *Nature*, 593(7858):249–254, May 2021.
- [4] Ran Wang, Xupeng Chen, Amirhossein Khalilian-Gourtani, Leyao Yu, Patricia Dugan, Daniel Friedman, Werner Doyle, Orrin Devinsky, Yao Wang, and Adeen Flinker. Distributed feedforward and feedback processing across perisylvian cortex supports human speech. December 2021.
- [5] Francis R. Willett, Erin M. Kunz, Chaofei Fan, Donald T. Avansino, Guy H. Wilson, Eun Young Choi, Foram Kamdar, Leigh R. Hochberg, Shaul Druckmann, Krishna V. Shenoy, and Jaimie M. Henderson. A high-performance speech neuroprosthesis. January 2023.
- [6] Sean L. Metzger, Kaylo T. Littlejohn, Alexander B. Silva, David A. Moses, Margaret P. Seaton, Ran Wang, Maximilian E. Dougherty, Jessie R. Liu, Peter Wu, Michael A. Berger, Inga Zhuravleva, Adelyn Tu-Chan, Karunesh Ganguly, Gopala K. Anumanchipalli, and Edward F. Chang. A high-performance neuroprosthesis for speech decoding and avatar control. *Nature*, 620(7976):1037–1046, August 2023.
- [7] Zhenhailong Wang and Heng Ji. Open vocabulary electroencephalography-to-text decoding and zero-shot sentiment classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5350–5358, 2022.
- [8] Yiqun Duan, Charles Zhou, Zhen Wang, Yu-Kai Wang, and Chin teng Lin. Dewave: Discrete encoding of eeg waves for eeg to text translation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [9] Debadatta Dash, Paul Ferrari, and Jun Wang. Decoding imagined and spoken phrases from non-invasive neural (meg) signals. *Frontiers in neuroscience*, 14:290, 2020.
- [10] Richard Csaky, Mats WJ van Es, Oiwi Parker Jones, and Mark Woolrich. Interpretable many-class decoding for meg. *NeuroImage*, 282:120396, 2023.
- [11] Gayane Ghazaryan, Marijn van Vliet, Aino Saranpää, Lotta Lammi, Tiina Lindh-Knuutila, Annika Hultén, Sasa Kivisaari, and Riitta Salmelin. Trials and tribulations when attempting to decode semantic representations from meg responses to written text. *Language, Cognition and Neuroscience*, pages 1–12, 2023.
- [12] Jo Hyejeong, Yang Yiqian, Juhyeok Han, Yiqun Duan, Hui Xiong, and Won Hee Lee. Are eeg-to-text models working? *arXiv preprint arXiv:2405.06459*, 2024.
- [13] Yiqian Yang, Yiqun Duan, Qiang Zhang, Renjing Xu, and Hui Xiong. Decode neural signal as speech. *arXiv preprint arXiv:2403.01748*, 2024.
- [14] Alexandre Défossez, Charlotte Caucheteux, Jérémy Rapin, Ori Kabeli, and Jean-Rémi King. Decoding speech perception from non-invasive brain recordings. *Nature Machine Intelligence*, 5(10):1097–1107, October 2023.
- [15] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR, 2023.
- [16] Laura Gwilliams, Graham Flick, Alec Marantz, Liina Pyllkkänen, David Poeppel, and Jean-Rémi King. Introducing meg-masc a high-quality magneto-encephalography dataset for evaluating natural speech processing. *Scientific Data*, 10(1), December 2023.
- [17] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [18] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [19] Sean L. Metzger, Jessie R. Liu, David A. Moses, Maximilian E. Dougherty, Margaret P. Seaton, Kaylo T. Littlejohn, Josh Chartier, Gopala K. Anumanchipalli, Adelyn Tu-Chan, Karunesh Ganguly, and Edward F. Chang. Generalizable spelling using a speech neuroprosthesis in an individual with severe limb and vocal paralysis. *Nature Communications*, 13(1), November 2022.
- [20] Yan Liu, Zehao Zhao, Minpeng Xu, Haiqing Yu, Yanming Zhu, Jie Zhang, Linghao Bu, Xiaoluo Zhang, Junfeng Lu, Yuanning Li, Dong Ming, and Jinsong Wu. Decoding and synthesizing tonal language speech from brain activity. *Science Advances*, 9(23), June 2023.

- [21] C Feng, L Cao, D Wu, E Zhang, T Wang, X Jiang, H Ding, C Zhou, J Chen, H Wu, et al. A high-performance brain-to-sentence decoder for logossyllabic language. 2023.
- [22] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, 2020.
- [23] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, 33:12449–12460, 2020.
- [24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [25] Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006.
- [26] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.10512*, 2023.
- [27] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- [28] Emilia P Martins and Theodore Garland Jr. Phylogenetic analyses of the correlated evolution of continuous characters: a simulation study. *Evolution*, 45(3):534–557, 1991.
- [29] Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Texygen: A benchmarking platform for text generation models. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 1097–1100, 2018.

Supplementary Material for MAD: Multi-alignment MEG-text decoding

A Dataset

The Gwilliams [16] dataset is described below:

A.1 Participants

- **Total Participants:** 27 English-speaking adults (15 females)
- **Age:** Mean = 24.8 years, SD = 6.4 years
- **Recruitment:** Subject pool of NYU Abu Dhabi
- **Consent and Compensation:** All provided written informed consent and were compensated
- **Health:** Reported normal hearing and no history of neurological disorders
- **Language:** All but one participant (S20) were native English speakers
- **Sessions:**
 - Majority (22 participants) performed two identical one-hour-long sessions
 - Sessions were separated by 1 day to 2 months
- **Ethics Approval:** Approved by the IRB ethics committee of NYU Abu Dhabi

A.2 Procedure

- **Recording Sessions:**
 - Duration: Each session lasted approximately 1 hour.
 - Equipment: Recorded with a 208 axial-gradiometer MEG scanner (Kanazawa Institute of Technology).
 - Sampling Rate: 1,000 Hz.
 - Filtering: Online band-pass filtered between 0.01 and 200 Hz.
 - Task: Participants listened to four distinct stories through binaural tube earphones (Aero Technologies) at a mean level of 70 dB sound pressure level.
- **Pre-Experiment Exposure:**
 - Participants were exposed to 20 seconds of each distinct speaker voice.
 - Purpose: To clarify session structure and familiarize participants with the voices.
- **Story Presentation Order:**
 - Assigned pseudo-randomly using a "Latin-square design."
 - Same order used for both recording sessions for each participant.
- **Attention Check:**
 - Participants answered a two-alternative forced-choice question every 3 minutes.
 - Example Question: "What precious material had Chuck found? Diamonds or Gold."
 - Average Accuracy: 98%, confirming engagement and comprehension.
- **MRI Scans:**
 - T1-weighted anatomical scans were performed after MEG recording if not already available.
 - Six participants did not return for their T1 scan.
- **Head Shape Digitization:**
 - Head shape digitized with a hand-held FastSCAN laser scanner (Polhemus).
 - Co-registered with five head-position coils.
 - Coil positions collected before and after each recording, stored in the 'marker' file.
 - Experimenter continuously monitored head position to minimize movement.

A.3 Stimuli

- **Stories:** Four English fictional stories selected from the Open American National Corpus:
 - ‘**Cable spool boy**’: 1,948 words, narrating two young brothers playing in the woods.
 - ‘**LW1**’: 861 words, narrating an alien spaceship trying to find its way home.
 - ‘**Black willow**’: 4,652 words, narrating the difficulties an author encounters during writing.
 - ‘**Easy money**’: 3,541 words, narrating two friends using a magical trick to make money.
- **Audio Tracks:**
 - Synthesized using Mac OS Mojave’s (c) text-to-speech.
 - Voices and speech rates varied every 5-20 sentences to decorrelate language from acoustic representations.
 - Voices used: ‘Ava’, ‘Samantha’, and ‘Allison’.
 - Speech rate: Between 145 and 205 words per minute.
 - Silence between sentences: Varied between 0 and 1,000 ms.
- **Story Segments:**
 - Each story divided into 5-minute sound files.
 - Random word list played approximately every 30 seconds, generated from unique content words of the preceding segment.
 - Very small fraction (<1%) of non-words introduced in natural sentences.
- **Task Definition:**
 - Each "task" corresponds to the concatenation of sentences and word lists.
 - All subjects listened to the same set of four tasks, in different block orders.

B Discussion about the main table

We used BLEU-1 [17], ROUGE-1-F [18], BertScore [27], CER [28], Self-BLEU [29] as metrics in the main table to show the capability of previous models and our models. However, as observed, NeuSpeech [13] model has the best score for ROUGE-1, Bert, CER, which is incredible, therefore we measured the Self-BLEU of this model, which is almost 100%, and found out NeuSpeech predicts almost the same sentence “He looked at me and said to me” all the time for different sentences in Supp. 1. Generation of this bad quality has best score on these three metrics, which means these three metrics are not effective in measuring the generation quality. Therefore, we think BLEU-1 the most reliable metric in this task for now. Besides, we randomly selected sentences, which is RS in the table, from the test set as another baseline, we found out that the BLEU-1 score is higher than NeuSpeech, which means the NeuSpeech model is not effective, which is very reasonable. After all, it seems that using BLEU score is the only reasonable metric of evaluating the quality of generated text.

As observed in the table, it is very clear that our MAD model is significantly higher than RS and NeuSpeech and Wav2vec2CTC on BLEU-1, which means our MAD model is effective on unseen text.

C More generated samples

We showed more generate samples here to show that we are not cherry-picking.

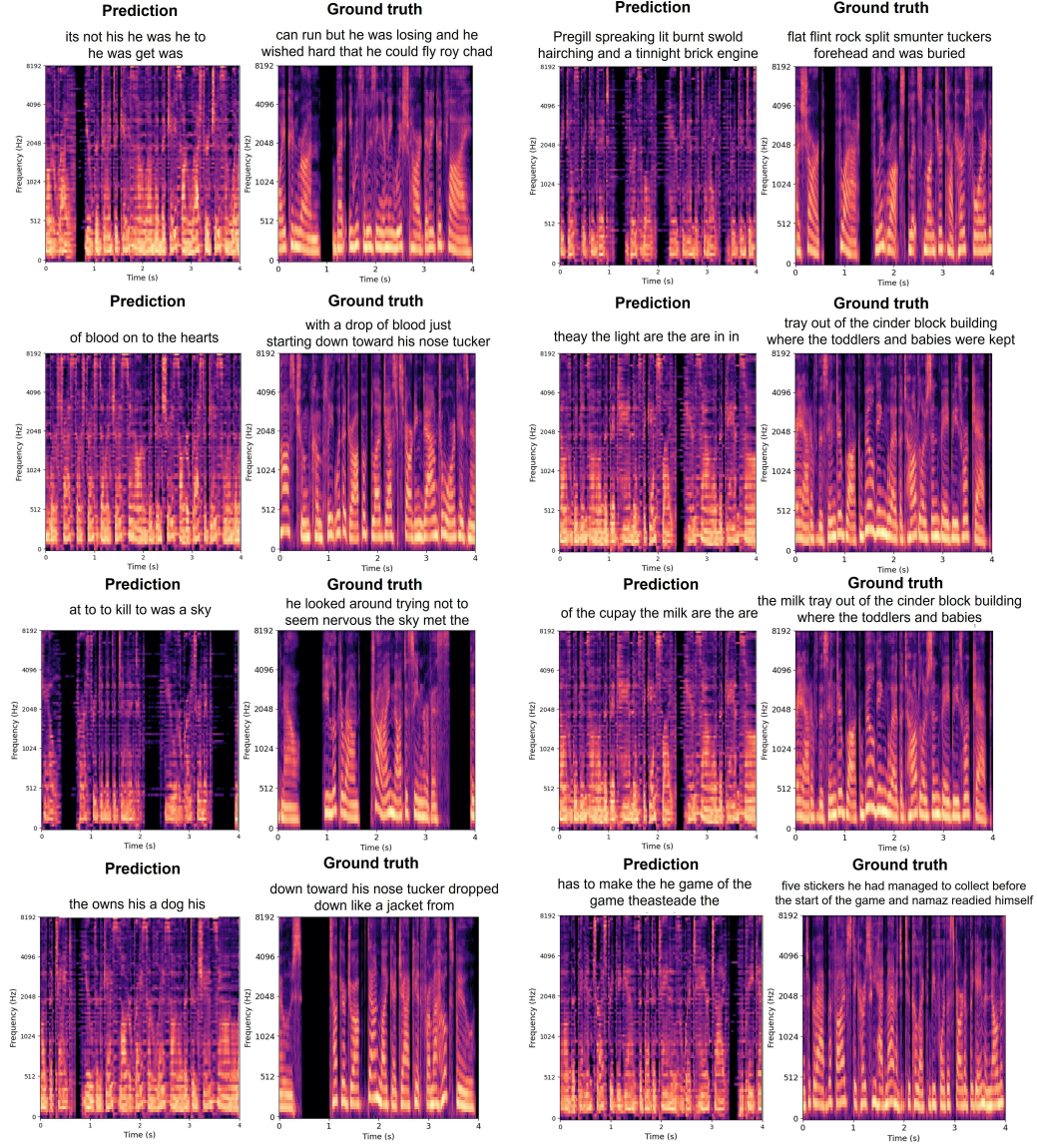


Figure 3. Eight sample examples of the test set. Prediction refers to Mel spectrograms generated by the brain module. Ground truth refers to Mel spectrograms of the audio signal processed by the whisper processor. The predicted text was generated using teacher forcing. These examples were produced using $L_m(mmd)$ with only a trainable brain module.

Listing 1. NeuSpeech [13] generation without teacher forcing.

```

1 start*****
2 Predicted: He looked at me and said to me,
3 True: were smelly thistles or cocklebur stems covered with spiked
4 end=====
5
6 start*****
7 Predicted: He looked at me and said to me,
8 True: or ordering Chad around or something. But since his fall the
9 year before,
10 end=====
11
12 start*****
13 Predicted: I'm not sure how to do it. It's just a little bit more
14 True: oldest boy in the playground, and the one who decided the rules
15 end=====
16
17 start*****
18 Predicted: He looked at me and said to me,
19 True: Spauw for fear of what was coming next. I'll make you fight.
20 Tucker
21 end=====
22
23 start*****
24 Predicted: he looked at me and said, I don't know what to do.
25 True: before, Roy had been shuffling and doing what he was told. Chad
26 end=====
27
28 start*****
29 Predicted: He looked at me and said to me,
30 True: for the tumbleweed to prove he wasn't a baby to Tucker. But as
31 much
32 end=====
33
34 start*****
35 Predicted: He looked at me and said to me,
36 True: walk really every something great blade over. Mama
37 end=====
38
39 start*****
40 Predicted: He looked at me and said to me,
41 True: other ready to step down into Chad's back. A sharp, Flat,
42 end=====
43
44 start*****
45 Predicted: He looked at me and said to me,
46 True: about gathering stickers himself. Roy was too
47 end=====
48
49 start*****
50 Predicted: He looked at me and said to me,
51 True: in shade and napped inside the walls. Then could wild and blink-
52 breath corner-hard
53 end=====

```


Listing 2. Wav2vec2CTC [14] generation.

```

1 start*****
2 Predicted: THLE'S HOAN BSFBHLAG'DS HON CITES HAG THOEANGLEN S QJRANGD
3 HOAND'S SORUESTHO E MRERLWOAINS HOAX TH
4 True: AND NAPPED INSIDE THE WALLS THEN COULD WILD AND BLINKBREATH
5 CORNERHARD
6 end=====
7 start*****
8 Predicted: SHROE BHOING TSEDTRAINS BBB
9 True: OF TIRES TWO BIG TRACTOR TIRES CAPPED OUT WITH ONE FROM A TRUCK
10 AND TWO SMALLER
11 end=====
12 start*****
13 Predicted: IES HO BHE HRORA SCIRCIND FBW
14 True: THAT OUT EITHER IT WAS ROY'S FAVORITE GAME NO
15 end=====
16 start*****
17 Predicted: AGSCHRONDSOUNE HIRS ON HOIN PHRORLI'S HEXSHIS B
18 True: ABOUT GATHERING STICKERS HIMSELF ROY WAS
19 end=====
20 start*****
21 Predicted: D JABWUISD BHOEND TE AUST THORE MLADS BHAXTS BMOIST OND F
22 True: TWO SMALLER ONES FROM CARS THE OLDER BOYS LAY AROUND IN
23 end=====
24 start*****
25 Predicted: CHORWALDES OE CSCRER BXSCOUE WONSTFBHE HOITS PR ENS
26 True: WASN'T CHICKEN YOU WANNA PLAY ROBOTS ROY ASKED CHAD
27 end=====
28 start*****
29 Predicted: BHI'S JMA
30 True: WHAT YOU SUCK CHAD SAID HE WISHED ROY
31 end=====
32 start*****
33 Predicted: SHOUDTIES BVIENT HOAS S
34 True: MAKING A DOOR TO THE SMALL ROOM INSIDE THE TALL TUMBLEWEED FLAG
35 end=====
36 start*****
37 Predicted: IDH HOASTD HIE' SCHORK SPHRERG 'S THOANS OABLWSDT'T XSCIED
38 HRIE HOER SPTHREALNINDSFOFTHES PHE CHOR HIER
39 True: WEAPONS ALLOWED ACCORDING TO HUMPTY DUMPTY NURSERY RULES
40 end=====
41 start*****
42 Predicted: SHOURX PHRERLNGDS FHOANS OMBLWSDT'T ESCED RIE HORN
43 SFTHRANINDSFOTS PHE CHOR CHIRE HINS HIND HOURXS TH
44 True: ALLOWED ACCORDING TO HUMPTY DUMPTY NURSERY RULES OF ENGAGEMENT
45 end=====
46
47
48
49
50

```

Listing 3. MAD generation with teacher forcing.

```

1 start*****
2 Predicted: orus said wast be a day but out
3 True: chad said he wished roy wouldnt fall for that gag every time get
4 end=====
5
6 start*****
7 Predicted: name is from his head on his head ofs
8 True: down his head rose and his eyes focused over chads shoulder out
9 roy
10 end=====
11
12 start*****
13 Predicted: be the smell times have at
14 True: until he could smell the dust several hated must staring brother
15 end=====
16
17 start*****
18 Predicted: he had not but though he was not a be down to
19 True: he wished he were there now even if he did have to sit next
20 end=====
21
22 start*****
23 Predicted: is sky the the ground side of the sky
24 True: the sky met the flat ground in all directions on the other side
25 of the chainlink fence
26 end=====
27
28 start*****
29 Predicted: the the up lift him know the the rest the that ist the fool
30 the he
31 True: to lift him and let him reach for the tumbleweed to prove he
32 wasnt a baby to tucker but
33 end=====
34
35 start*****
36 Predicted: sound is the mouth ist been
37 True: a sick sound but the thing in his head hadnt worked
38 end=====
39
40 start*****
41 Predicted: the of the top a red medal
42 True: out of the top of the black fort like a gold headed monster
43 end=====
44
45 start*****
46 Predicted: the roy him he name was be
47 True: out after him roy chad called but his voice would
48 end=====
49
50 start*****
51 Predicted: soldiers astronautss a on be us and
52 True: for soldiers and astronauts and its vote going to help roy
53 end=====

```

D Ethics

D.1 Safety

Our MEG-to-text translation technology is designed to assist individuals with severe speech and motor impairments by translating brain signals into text. While this technology has the potential to greatly improve quality of life, we acknowledge the possibility of misuse. However, there are no foreseeable situations where the direct application of our technology could harm, injure, or kill people. We do not develop or intend to develop applications that increase the lethality of weapons systems.

D.2 Security

We recognize the importance of securing brain-computer interface systems. Future research should include thorough risk assessments to identify and mitigate potential security vulnerabilities. We recommend employing robust encryption methods and secure data transmission protocols to protect against unauthorized access and ensure the safety of users' neural data.

D.3 Discrimination

Our technology aims to provide equal accessibility to communication for individuals with speech and motor impairments. We are committed to ensuring that our MEG-to-text translation system does not facilitate discrimination or exclusion. We will continuously monitor and test our models to prevent biases that could negatively impact service provision in healthcare, education, or financial sectors.

D.4 Surveillance

We adhere strictly to local laws and ethical guidelines regarding data collection and analysis. Our research does not involve bulk surveillance data. We obtained data from public dataset [16], and we do not predict protected categories or use data in ways that endanger individual well-being.

D.5 Deception & Harassment

Our technology is designed with safeguards to prevent its misuse in deceptive or harmful interactions. We implement verification mechanisms to detect and prevent impersonation and fraudulent activities. We actively work to ensure our system cannot be used to promote hate speech, abuse, or influence political processes maliciously.

D.6 Environment

We are mindful of the environmental impact of our research. While our work primarily involves computational resources, we strive to optimize our algorithms to minimize energy consumption. We do not engage in activities that promote fossil fuel extraction or increase societal consumption. Our focus is on developing sustainable and efficient technologies.

D.7 Human Rights

Our research adheres to ethical standards and legal requirements, ensuring that it does not facilitate illegal activities or deny individuals their rights to privacy, speech, health, liberty, security, legal personhood, or freedom of conscience or religion. We are committed to protecting and promoting human rights through our work.

D.8 Bias and Fairness

Our goal is to create fair and inclusive technology that benefits all users equally, regardless of their background.

E Broader Impacts

E.1 Positive Impacts

Improved Communication for Individuals with Disabilities Our MEG-to-text translation technology has the potential to significantly improve the quality of life for individuals with severe speech and motor impairments. By enabling these individuals to communicate effectively, we can help them achieve greater independence, participate more fully in society, and reduce their reliance on caregivers.

Advancements in Neurotechnology This research contributes to the broader field of neurotechnology, advancing our understanding of brain activity and its relationship to language. These advancements could lead to new therapies and interventions for a variety of neurological conditions, potentially benefiting a wide range of patients.

Innovation and Economic Growth The development and commercialization of advanced neurotechnologies can stimulate economic growth by creating new industries and job opportunities. This innovation can drive progress in related fields such as healthcare, robotics, and artificial intelligence, fostering a collaborative and dynamic technological ecosystem.

E.2 Negative Impacts

Privacy and Security Risks The collection and analysis of neural data pose significant privacy and security concerns. Unauthorized access to such sensitive information could lead to misuse, including identity theft or unauthorized surveillance. Ensuring robust data protection measures is crucial to mitigate these risks.

Potential for Misuse There is a risk that the technology could be misused for purposes other than its intended therapeutic applications. For instance, it could be exploited for invasive surveillance or to manipulate individuals by decoding their thoughts without consent. Strict ethical guidelines and regulations are necessary to prevent such misuse.

Bias and Discrimination If not carefully managed, the technology could inadvertently encode or exacerbate existing biases. For example, if the training data is not representative of diverse populations, the system might perform poorly for certain groups, leading to unequal access to its benefits. Ongoing efforts to ensure fairness and inclusivity are essential to address these concerns.