

DeWave: Discrete EEG Waves Encoding for Brain Dynamics to Text Translation

Yiqun Duan^{1*}, Jinzhao Zhou¹, Zhen Wang², Yu-Kai Wang¹, Chin-Teng Lin^{1†}

¹GrapheneX-UTS HAI Centre, Australian Artificial Intelligence Institute,

Faculty of Engineering and Information Technology

University of Technology Sydney, Ultimo, NSW 2007

²School of Computer Science, The University of Sydney, Camperdown NSW 2050

{yiqun.duan, jinzhao.zhou}@student.uts.edu.au, zwan4121@uni.sydney.edu.au

yukai.wang@uts.edu.au, chin-teng.Lin@uts.edu.au

Abstract

The translation of brain dynamics into natural language is pivotal for brain-computer interfaces (BCIs). With the swift advancement of large language models, such as ChatGPT, the need to bridge the gap between the brain and languages becomes increasingly pressing. Current methods, however, require eye-tracking fixations or event markers to segment brain dynamics into word-level features, which can restrict the practical application of these systems. To tackle these issues, we introduce a novel framework, DeWave, that integrates discrete encoding sequences into open-vocabulary EEG-to-text translation tasks. DeWave uses a quantized variational encoder to derive discrete codex encoding and align it with pre-trained language models. This discrete codex representation brings forth two advantages: 1) it realizes translation on raw waves without marker by introducing text-EEG contrastive alignment training, and 2) it alleviates the interference caused by individual differences in EEG waves through an invariant discrete codex with or without markers. Our model surpasses the previous baseline (40.1 and 31.7) by 3.06% and 6.34%, respectively, achieving 41.35 BLEU-1 and 33.71 Rouge-F on the ZuCo Dataset. This work is the first to facilitate the translation of entire EEG signal periods without word-level order markers (e.g., eye fixations), scoring 20.5 BLEU-1 and 29.5 Rouge-1 on the ZuCo Dataset.

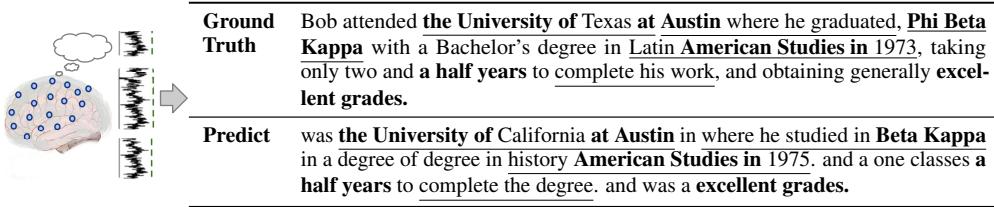


Figure 1: Overall illustration of translating EEG waves into text ²through quantised encoding.

1 Introduction

Decoding brain states into comprehensible representations has long been a focal point of research [20, 38, 8, 54]. Electroencephalogram (EEG) signals are particularly favored by researchers due to their non-invasive nature and ease of recording [31, 46]. Traditional EEG decoding techniques

*is the first author, † is the corresponding author

²This visualization still keeps teacher-forcing evaluation for a fair comparison with previous methods.

largely focus on classifying brain states into restricted categories like Motor Imaginary (MI) [37, 7], Emotion [19, 47], Robotic Control [44, 55], and Gaming [28, 24]. However, these labels, bound to specific tasks, are insufficient for broad-based brain-computer communication. Consequently, there has been a surge of interest in brain-to-text (speech) translation in recent years. As the current trend leans towards large language models (LLM) [4, 12, 30] showcasing increasingly generalized intelligence capabilities, it becomes crucial to delve into ways of bridging the gap between brain signals and natural language representation. However, this area remains under-explored.

The early work in brain-to-text translation [14, 1, 26, 41] relied on external event markers like handwriting or eye-tracking fixations to segment whole brain signals into fragmentary features. This methodology treated the task as word-level classification on a small, closed vocabulary set, with each time step analyzed individually. Both invasive [14, 42] (ECoG) and non-invasive [33] (EEG) brain signals have been used in these approaches. Notably, utilizing handwriting as event markers on invasive signals has led researchers [49] to achieve state-of-the-art (SOTA) recognition accuracy on closed-set character-level recognition. Wang [48] expanded the vocabulary size substantially and demonstrated the feasibility of open-vocabulary brain-to-text translation by employing pre-trained language models with word-level EEG features. However, limitations persist. The order of event markers, particularly eye-tracking fixations used to segment EEG waves into word-level features, may not match the natural word order in language³. Moreover, current methods do not have the capacity for direct text translation.

In this paper, we present Discrete EEG Waves Encoding for Brain Dynamics to Text Translation (DeWave), a pioneering framework depicted in Figure 2. DeWave uses a vector quantized variational encoder to transform EEG waves into a discrete codex, linking EEG waves to tokens based on their proximity to codex book entries. This method offers two key advantages: it addresses significant distribution variances in EEG waves across individuals [29, 40, 9], and rectifies order mismatches between raw wave sequences and text without eye-tracking markers. Our raw waves encoder is guided by both self-reconstruction and a contrastive supervision alignment between text embeddings and vectorized raw waves. To navigate the challenges of training with limited parallel data, DeWave leverages large-scale pre-trained language models [6, 36, 3], specifically employing BART [21], which combines BERT’s bidirectional context with GPT’s left-to-right decoder. Notably, our discrete codex aligns more closely with actual language tokens than continuous EEG features, serving as an interpretable bridge between EEG input and the language model.

Experiments employ non-invasive EEG signals and data from the ZuCo dataset [16], a large-scale public resource that records eye-tracking and EEG during natural reading tasks. Notably, DeWave can be generalized for both word-level EEG features and raw EEG wave translation. With a robust codex representation, this work pioneers in translating the entire time period of EEG signals without the need for word-level order markers such as eye fixations. We assess the performance using standard translation metrics [34, 25]. With word-level EEG features, DeWave attains 42.8 BLEU-1 and 34.9 Rouge-1, surpassing the previous baseline (40.1 and 31.7) by 6.73% and 10.09% respectively on the ZuCo Dataset. For raw EEG waves without event markers, DeWave achieves 20.5 BLEU-1 and 29.5 Rouge-1. This work’s contributions can be summarized in three main points.

- This paper introduces discrete codex encoding to EEG waves and proposes a new framework, DeWave, for open vocabulary EEG-to-Text translation.
- By utilizing discrete codex, DeWave is the first work to realize the raw EEG wave-to-text translation, where a self-supervised wave encoding model and contrastive learning-based EEG-to-text alignment are introduced to improve the coding ability.
- Experimental results suggest the DeWave reaches SOTA performance on EEG translation, where it achieves 41.35 BLUE-1 and 33.71 Rouge-1, which outperforms the previous baselines by 3.06% and 6.34% respectively.

2 Related Works

The key to decoding natural language from EEG signals is good representations. Existing work for EEG-to-text representation can be categorized into hand-crafted representations and deep-learning

³For example, the ZuCo Dataset collects data by simultaneously recording eye-tracking fixation and brain waves during reading tasks. However, the order of eye-fixations and spoken words may not always coincide.

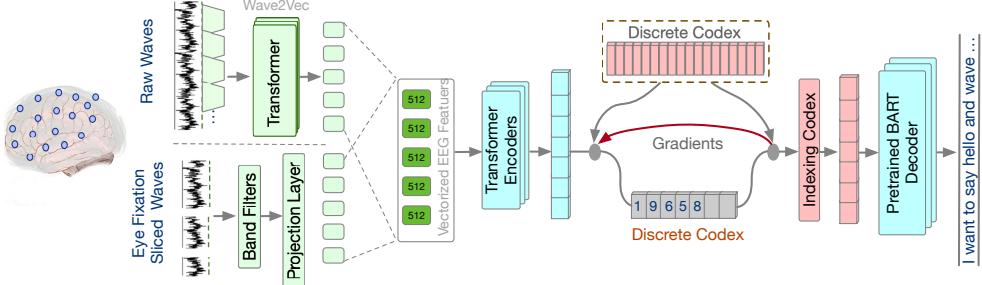


Figure 2: The DeWave model structure involves vectorizing either word-level EEG features or raw EEG waves into embeddings (Section 3.3). The vectorized features are then encoded into a latent variable $\mathbf{z}_c(\mathcal{X})$, which is converted into a discrete latent $\mathbf{z}_q(\mathcal{X})$ through codex indexing. Finally, a pre-trained BART model translates this discrete codex representation into texts.

representations. Earlier works make use of traditional methods to obtain hand-crafted representations from a sequence of EEG signals. A continuous feature representation is extracted using methods such as statistical features [53], correlation coefficient matrix from fast Fourier transform (FFT) components [43, 51, 50], wavelet transform features [18, 32, 54], and Mel-frequency cepstral coefficients (MFCCs) [5]. These hand-crafted representations can be used to establish a mapping between EEG segments and words using distance-based methods such as linear discriminative analysis (LDA) or Support Vector Machine (SVM).

On the other hand, deep learning methods, such as [52], utilize a convolutional neural network (CNN) and a Long Short-Term Memory (LSTM) to learn deep features and perform classification of user’s instructions. However, the direct representation is often insufficient to discriminate a larger number of imagined text categories. In [39], EEG signals are first classified into six phonological categories as the intermediate state using a deep convolutional autoencoder. Then latent features are used as input for another encoder to predict a total of 11 speech tokens. Recently, [48] proposed a framework that uses a multi-layer transformer encoder to project word-level EEG feature sequences into EEG embeddings. Then a pre-trained BART model is used to decode these embeddings into words.

3 Method

The overall process of DeWave is illustrated in Figure 2, where the word-level or raw EEG features are vectorized into sequences embedding and are fed into the discrete codex. The language model generates translation output based on discrete codex representation.

3.1 Task Definition

Given a sequence of word-level EEG features \mathcal{E} , the aim is to decode the corresponding open-vocabulary text tokens \mathcal{W} . These EEG-Text pairs $\langle \mathcal{E}, \mathcal{W} \rangle$ are collected during natural reading, as defined in Section 4.1. We delve into two task settings: (1) Word-level EEG-to-Text Translation, where EEG feature sequences \mathcal{E} are fragmented and re-ranked based on eye-fixation \mathcal{F} aligned with each word token w in sequence \mathcal{W} ; and (2) Raw EEG Waves to Text Translation, where EEG feature sequences \mathcal{E} are directly vectorized into embedding sequences for translation without any event markers, a more challenging but practical real-time setting. DeWave is the pioneering work in this latter task.

3.2 Discrete Codex

Discrete representation is first proposed in VQ-VAE [45]. DeWave is the first work to introduce discrete encoding into EEG signal representation. The discrete representation could benefit both the word-level EEG features and the raw EEG wave translation. Introducing discrete encoding into brain waves could bring two aspects of advantages. 1) It is widely accepted that EEG features have a strong data distribution variance across different human subjects. Meanwhile, the datasets can only have samples from a few human subjects due to the expense of data collection. This severely weakened the generalized ability of EEG-based deep learning models. By introducing discrete encoding, we

could alleviate the input variance to a large degree as the encoding is based on checking the nearest neighbor in the codex book. 2) The codex contains fewer time-wise properties which could alleviate the order mismatch between event markers (eye fixations) and language outputs. Meanwhile, because of this property of codex represents, DeWave is the first work that could realize the direct translation from raw EEG waves without any event marker to give the order.

Inference: Given the EEG waves \mathcal{E} , it is first vectorized into embedding as introduced in Section 3.3 with $(\mathcal{X} = \Theta(\mathcal{E}, \mathcal{F}))$ or without $(\mathcal{X} = \Theta(\mathcal{E}))$ eye fixations \mathcal{F} , where \mathcal{X} is the embedding sequence. A codex book $\{\mathbf{c}_i\} \in \mathbb{R}^{k \times m}$ is initialized with number k of latent embedding with size m . The vectorized feature \mathcal{X} is encoded into $\mathbf{z}_c(\mathcal{X})$ through a transformer encoder. The discrete representation is acquired by calculating the nearest embedding in the codex of input embedding $\mathbf{x} \in \mathcal{X}$ as shown in Equation 1.

$$\mathbf{z}_q(\mathcal{X}) = \{\mathbf{z}_q(\mathbf{x})\}_i, \quad \mathbf{z}_q(\mathbf{x}) = \mathbf{c}_k, \quad k = \operatorname{argmin}_j \|\mathbf{z}_c(\mathbf{x}) - \mathbf{c}_j\|_2 \quad (1)$$

Different from the original VQ-VAE which decodes the original input, Dewave directly decodes the translation output given the representation $\mathbf{z}_q(\mathcal{X})$. Given a pre-trained language model, the decoder predicts text output with $P(\mathcal{W}|\mathbf{z}_q(\mathcal{X}))$.

Learn: The codex is like a bridge connecting the vectorized EEG feature and the language model. Compared to learning direct EEG-to-text relation, DeWave learns a better discrete codex for the language model. It is easier to learn since We learn the discrete codex by the combination of the loss functions in three parts,

$$L = -\log(p(\mathcal{W}|\mathbf{z}_q(\mathcal{X}))) + \|\mathbf{sg}[\mathbf{z}_c(\mathbf{x})] - \mathbf{z}_q(\mathbf{x})\|_2^2 + \beta \|\mathbf{z}_c(\mathbf{x}) - \mathbf{sg}[\mathbf{z}_q(\mathbf{x})]\|_2^2 \quad (2)$$

where the loss maximize the log-likelihood of language outputs $\log(P(\mathcal{W}|\mathbf{z}_q(\mathcal{X})))$ and minimize the distance between latent variable \mathbf{z} and the codex value \mathbf{z} . Here, the \mathbf{sg} denotes the stop gradients. The learning is robust for β from 0.1-2.0, where we set it as 0.2 throughout the training process.

3.3 EEG Vectorization

Word-Level EEG Features With Event Markers: The EEG waves are first sliced into fragments according to the eye-tracking fixation of word sequences given in the annotation. Similar to [48], we calculate the statistical result of four frequency band filters, Theta band (5-7Hz), the Alpha band (8-13Hz), the Beta band (12-30Hz), and Gamma band (30Hz-) [27] to get the statistic frequency features of each fragment. It is noted that although different fragments may have different EEG window sizes, the statistical results are the same (embedding size 840). A multi-head transformer layer is applied to project the embedding into feature sequences with latent size 512.

Raw EEG Waves: Self-Guided Waves to Discrete Codex Our self-supervised EEG wave encoder transforms raw EEG signals into a sequence of embeddings [2, 56] as illustrated in Figure 3. It has two guiding principles: Self-Reconstruction, where the encoder is trained to transform and subsequently reconstruct the original waveforms from discrete codices; and Text Alignment, where the codices' encoding is semantically aligned with word vectors, fostering the development of text-aligned EEG signal representations.

For structure-wise, a conformer-based multi-layer encoder with specially designed hyperparameters is employed. The one-dimensional convolution layer processes the EEG waves to generate the embedding sequence ⁴, fusing the EEG channels into a unique embedding for each period. We apply bi-directional transformer attention layers to the sequence to capture temporal relations.

For the reconstruction process, a decoder transformer and transpose convolution structure then convert these discrete embeddings back into raw waves. Given the reconstruction process as $\tilde{\mathcal{X}} = \phi(\mathbf{z}_q(\mathcal{X}))$, the self-supervised loss could be modified into:

$$L_{\text{wave}} = \frac{1}{n} \sum_n (\phi(\mathbf{z}_q(\mathcal{X})) - \mathcal{X})_n^2 + \|\mathbf{sg}[\mathbf{z}_c(\mathbf{x})] - \mathbf{z}_q(\mathbf{x})\|_2^2 + \beta \|\mathbf{z}_c(\mathbf{x}) - \mathbf{sg}[\mathbf{z}_q(\mathbf{x})]\|_2^2, \quad (3)$$

where the model calculates the mean square loss between the reconstructed wave and the wave ground truth to perform self-supervised training.

To obtain a semantically coherent codex, we introduce a cross-modality contrastive learning approach distinct from prevalent methods. Unlike CLIP [35, 23, 22] that contrasts CLS embeddings between

⁴perception field is roughly 200ms with overlap 100ms for each embedding

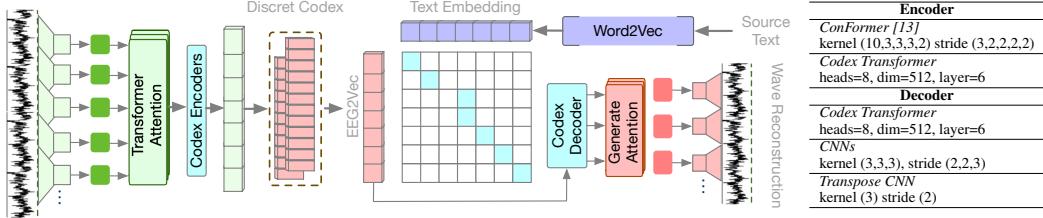


Figure 3: The image demonstrates the process of self-supervised pre-training for raw waves. The left subgraph details our strategy for directing the encoder, utilizing both self-reconstruction and text alignment through contrastive learning.

sample pairs within a mini-batch, our approach operates within a single EEG-text pair. We contrast the EEG-tokenized codex embeddings sequence z_q with the text embeddings sequence z_t . Assuming the raw wave feature extractors can produce a token sequence in an **organized chronological order**, we treat the diagonal EEG codex and text word2vec encoding pairs as positive pairs within the sequences. All other pairs are considered negative. The model is trained to minimize the distance between embeddings of positive pairs and maximize that of the negative pairs. For a given EEG-text pair (i, j) , the loss is defined in Equation 4

$$L_{\text{contrast}} = -\frac{1}{n} \sum \log \left[\frac{\exp(\mathbf{s}_{ii}/\tau)}{\sum_{k=1}^N \exp(\mathbf{s}_{ik}/\tau)} \right], \mathbf{s}_{ij} = \mathbf{z}_q(\mathbf{x}^i)^T \mathbf{z}_t(j) \quad (4)$$

Here, τ is the temperature parameter, N is the total number of EEG-text pairs in a batch, and the sum in the denominator is over all N EEG-text pairs and N text-EEG pairs. The EEG embeddings are expected to correctly match with their corresponding text while distinguishing them from mismatched EEG-text pairs. The total loss then becomes a combination of the original Wave2Vec loss and the contrastive loss as $L_{\text{total}} = L_{\text{wave}} + \alpha L_{\text{contrast}}$.

By this means, the model not only learns to reconstruct the EEG signal but also learns a robust representation of the signal that aligns with the corresponding text embeddings. This cross-modal learning can potentially improve the translation system by bridging the gap between EEG signals and the semantic content of the text.

3.4 Language Model

We used large-scale text corpus pre-trained BART [21] as the generative language model for translation output. As the EEG-to-text translation data is quite limited, leveraging the BART model could introduce prior knowledge of text relations. In that case, the translation system only needs to learn a codex representation for the language model, which is easier to learn. The codex representations are fed into pre-trained⁵ BART model and get the output hidden states. A fully connected layer is applied on the hidden states to generate English tokens from pre-trained BART vocabulary \mathcal{V} .

3.5 Training Paradigm

DeWave is trained through a multi-stage process, where the training process is illustrated in Appendix C. In the first stage, we do not involve the language model in weight updates. The target of the first stage is to train a proper encoder projection θ_{codex} and a discrete codex representation \mathcal{C} for the language model. In the second stage, the gradient of all weights, including language model θ_{BART} is opened to fine-tune the whole system.

4 Experiments

4.1 Dataset

DeWave utilize both ZuCo 1.0 [15] and 2.0 [17] for experiments. The dataset simultaneously recorded the text and EEG corpus during Normal Reading (NR) and Task-Specific Reading (TSR) tasks. The

⁵<https://huggingface.co/facebook/bart-large>

EEG waves are collected with a 128-channel system under a sampling rate 500Hz through a frequency band filter from 0.1Hz to 100Hz. However, after the noise canceling process, only 105 channels [15] are used for translation. Similar to [48], we slice the EEG wave according to the eye fixation and calculate the frequency features. For raw EEG waves, the signal is normalized into a value range of 0-1 for decoding. The reading task’s data are divided into the train (80%), development (10%), and test (10%) respectively by 10874, 1387, and 1387 unique sentences with no intersections. Please refer to Appendix A for a detailed description.

4.2 Implementation Details

For word-level EEG features, we use the 56 tokens each with an 840 embedding size. For raw EEG waves, we clip or pad the EEG waves up to sample point 5500 with a constant value of zero. A transformer layer with head number 8 and a 1×1 convolutional layer are combined to fuse multiple EEG channels into an embedding sequence with size 512. DeWave uses a codex with size 2048 where each codex latent is an embedding with size 512. The ablation study (Section 4.6) gives a discussion about the codex size. All models are trained on Nvidia V100 and A100 GPUs. For the self-supervised decoding for raw waves, we use a learning rate of 5e-4 and a VQ coefficient of 0.25 for training 35 epochs. For training the codex (stage 1), DeWave uses a learning rate of 5e-4 for 35 epochs. For finetuning the translation (stage 2), DeWave uses a learning rate of 5e-6 for 30 epochs. We use the SGD as the optimizer for training all the models. Due to limited space, refer to Appendix B for more details.

4.3 Evaluation Metrics

We evaluate translation performance using NLP metrics, BLEU and ROUGE, as shown in Table 1. For word-level EEG features, we compare our results to EEG-to-Text [48], maintaining a consistent language model for fairness. In the absence of methods for raw EEG waves, we establish a baseline EEG-to-Text † by segmenting the entire EEG waves into a sequence embedding using a 200ms time window with a 100ms overlap. We adapt Wave2Vec, originally developed for speech recognition, to brain waves and compare it with our approach, DeWave. Furthermore, we adapt unsupervised raw EEG waves classification methods BENDR [11] and SCL [10], using SSL pre-training and feature extraction for comparison, underscoring the impact of discrete encoding.

Table 1: Evaluation metrics of EEG-to-Text translation under both word-level features input and raw waves input, where +Contrastive denotes boost encoder with $L_{contrast}$ mentioned in Sec. 3.3. For a fair comparison, these results keep the same teacher-forcing evaluation setting as EEG-to-Text [48].

Source	Method	BLEU-N (%)				ROUGE-1 (%)		
		N=1	N=2	N=3	N=4	R	P	F
Word-level features	EEG-to-Text [48]	40.12	23.18	12.61	6.80	28.84	31.69	30.10
	DeWave	41.35	24.15	13.92	8.22	28.82	33.71	30.69
	EEG-to-Text † [48]	13.07	5.78	2.55	1.10	15.22	18.08	16.36
	Wave2Vec [2]	18.15	8.94	3.89	2.04	18.96	23.86	20.07
	BENDR [11]	18.48	9.16	4.05	2.15	19.03	25.22	21.18
	DeWave	20.51	10.18	5.16	2.52	21.18	29.42	24.27
Raw waves	DeWave+Contrastive	21.09	10.69	5.88	3.04	22.01	29.95	24.68

Word-Level EEG Features: For the word-level EEG feature, we observe that introducing the discrete brain representation could help DeWave reach BLEU-{1, 2, 3, 4} scores of 41.35, 24.15, 13.92, and 8.22, which respectively outperform the previous baseline by 1.23 (+3.06%), 0.97 (+4.18%), 1.31 (+10.38%) and 1.54 (+18.73%). It is observed that the increasing ratio is more significant for larger grams evaluation. DeWave achieve ROUGE-1 score 28.84 (R), 31.69 (P), and 30.10 (F) which outperforms the previous baseline by -0.02 (-0.06%), 1.98 (+6.35%), and 0.59 (+1.9%). The increase of the matrix is higher for longer phrases (3-gram and 4-gram). This is rational since the single-word level EEG features are sliced by eye fixation during reading, which may naturally contain noises since the human subject may not really think about the corresponding words when looking at them.

Raw Waves: DeWave represents an unprecedented endeavor in translating raw EEG waves directly, obviating the need for event markers. This is achieved through the application of a learned Wave2Vec model, which converts waves into discrete codex tokens. When benchmarked against

the baseline method - which merely slices the wave according to a time window - DeWave exhibits marked improvement. It attains BLEU-1, 2, 3, 4 scores of 20.51, 10.18, 5.16, and 2.56, respectively, surpassing the baseline by margins of 7.44 (+56.92%), 4.44 (+76.1%), 2.61 (+102.35%), and 1.42 (+129.09%). This underscores the superiority of a learned discrete embedding representation over the rudimentary time-window baseline. In comparison with cutting-edge self-supervised learning (SSL) methods, DeWave consistently outperforms them, including contrastive learning (CL)-based methods such as SCL [10], the original Wave2Vec [2], and BENDR [11].

4.4 Cross-Subject Performance

Cross-subject performance is of vital importance for practical usage. To further illustrate the performance variance on different subjects, we train the model by only using the data from subject YAG and test the metrics on all other subjects. The results are illustrated in Figure 4, where the radar chart denotes the performance is stable across different subjects. Please refer to Appendix G for more detailed results.

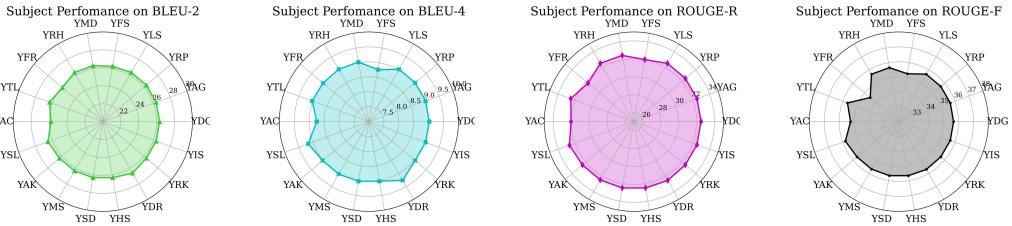


Figure 4: The cross-subjects performance on ZuCo dataset.

4.5 Generated Samples

In Table 2, we display visualized examples of text generated from unseen EEG signals. Despite the challenge of thought translation and limited prior research, our model yields meaningful results, aligning keywords and forming similar sentence structures, although it may not yet match traditional language translation tasks.

The model is more adept at matching verbs than nouns. For instance, **implies**” vs. **implies**”, **the author who wrote it**” vs. **the man who wrote it**” both effectively convey the intended sentiment of the sentence. However, when it comes to nouns, we observe a tendency towards synonymous pairs rather than precise translations, such as **the man**” vs. **the author**”, **Burroughs**” vs. **Heroughs**”, **edition**” vs. **version**”. Our analysis suggests two potential causes for this. First, when the brain processes these words, semantically similar words might produce similar brain wave patterns. Given the inherent noise in brain waves, the codex might group these features under the same value. Second, the volume of EEG-to-Text pairs available for training is significantly smaller than that for traditional language translation tasks. Hence, some degree of error in translating unseen nouns or sentences is to be expected.

Translation on raw EEG waves is naturally harder than word-level translation, as it lacks eye fixation to suggest the relationship between the period of waves and the word target. Table 2 meet our expectation that the results on raw EEG waves are not as good as those on word-level features, especially on the real semantic meaning of the sentence (sample (2) and sample (5) have the same target for comparison). However, the translation still could output the correct translation of certain words, such as **“much of”** vs. **“much of”**, **“individual”** vs. **“individual”**, and **“more complicated story”** vs. **“more exciting thing”**. Although EEG-to-text is a hard topic, DeWave suggests the feasibility of translation improvements.

4.6 Ablation Study

Discrete Codex DeWave encodes EEG waves into a discrete codex, aiming for a language model-friendly representation. We evaluated performance against varying codex sizes (1024 to 8192) to ascertain if larger sizes yield better results. As depicted in Figure 5, we found no strong correlation between codex size and model performance. A codex size of 2048 yielded the highest BLEU score average, and while the ROUGE score slightly improved with larger sizes, there was no clear evidence that increasing codex size consistently enhanced performance.

Table 2: Translation results on the unseen EEG waves, where **bold** denotes a correct match between ground truth and our prediction. Underline denotes a fuzzy match with similar semantic meanings. For a fair comparison, these results keep the same teacher-forcing evaluation setting as EEG-to-Text [48]. This means the decoding process eliminates accumulated errors and just predicts the current token with the GT token from the last step.

Decoding Results with Eye Fixation Assistance	
(1)	Ground Truth: Everything its title implies , a standard-issue crime drama spat out from the <u>Tinseltown</u> assembly line... Prediction: is own implies , including great of <u>issue</u> , novel of the beginning <u>towns</u> ...
(2)	Ground Truth: “The Kid Stays in the Picture” is a great story, terrifically told by the man who wrote it but this Cliff Notes <u>edition</u> is a cheat . Prediction: The film “says in the Game” is a film about but movie was written, the author who wrote it. also its <u>version</u> is a cheat .
(3)	Ground Truth: During Kerouac’s time at Columbia University, Burroughs and Kerouac got into trouble with the law for failing to report <u>a murder</u> ; this incident <u>formed</u> the basis of a mystery novel ... Prediction: Kerouac’s time at the, Heroughs and Kerouac were along a for the police for their to pay the <u>murder</u> . they <u>led</u> the basis of the lawsuit ...
Decoding Results with Raw Waves	
(4)	Ground Truth: Every individual will see the movie through the prism of <u>his or her own beliefs</u> and prejudices ... Prediction: Everyday individual is their results. their eyes of <u>his or her own personal</u> . desires. and ...
(5)	Ground Truth: “The Kid Stays in the Picture” is a great story, terrifically told by the man who wrote it but this Cliff Notes <u>edition</u> is a cheat . Model Output: The Price’s says in the Middle. and a common deal. and for good. this moment. made it is <u>still</u> little.
(6)	Ground Truth: much of this well-acted but dangerously slow thriller feels like a preamble to a bigger, more complicated story, Model Output: much of this is-being but not over. like it disaster-ble to a new and more exciting thing.

Raw EEG wave performance fluctuates noticeably with codex size. Performance improves when increasing the codex size from 1024 to 2048, but any larger size reduces performance. This variation may result from the training formation; word-level EEG data, selected by eye fixations, contains less noise than raw waves. We hypothesize that our current training data may be insufficient for larger codex sizes. Additionally, our experiments indicate that the latent dimension of the codex doesn’t significantly affect performance.

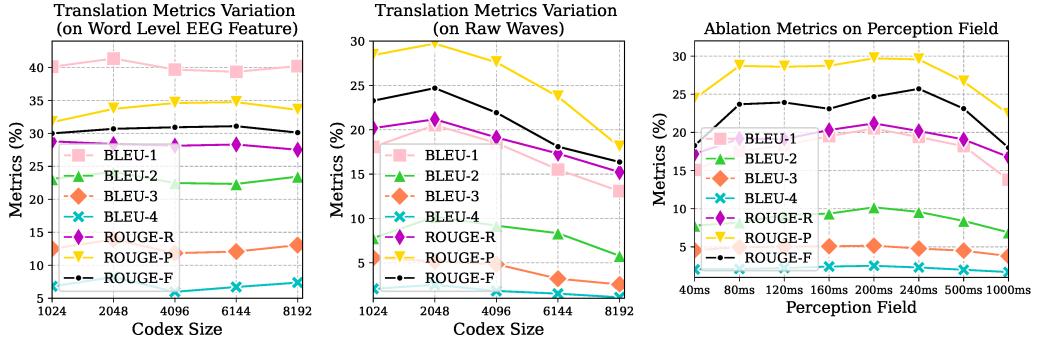


Figure 5: Ablation study on different codex sizes and perception fields (raw waves).

Perception Time Window We also conduct the ablation study on the model structure for the wave2vec model illustrated in Figure 5. As the model utilize a multi-layer CNN model to slide through the raw waves, the model compresses the waves for perception. The compress ratio decided how large the perception field is for each extracted embedding feature. As described in Section 3.3, the model utilized a perception field of 200ms with an overlap of 100ms. We conduct an ablation study of different perception fields and report it in Figure 5. Where it is observed that the model performance is significantly lower when the perception field is smaller than 80ms or larger than 240ms. The model could achieve similar results in the perception field 120ms to 240ms. The model reaches a small peak around 200ms to 240ms. We think this phenomenon is rational since the normal reading speed for humans is around 160-400 words per minute (reading speed may vary from different

material and human subjects). In other words, the reading period for each word is 150-375ms on average, which roughly meets our observation between 200ms-240ms.

Self-Supervision Initialization Theoretically, we could pre-train the codex by introducing a decoder and calculate the reconstruction loss with the original input for both word-level EEG and raw EEG waves. We prefix other parameters as our best setting with codex size 2048 and compare the impact in Table 3.

Table 3: Ablation on self-supervised pre-trained codex weights.

Pretrain	BLEU (%)			ROUGE-1 (%)		
	N=1	N=2	N=3	R	P	F
Word-level features	✓	41.35	24.15	13.92	28.82	33.71
	✗	40.71	22.94	12.40	28.01	34.08
Raw waves	✓	20.51	10.18	5.16	21.18	29.42
	✗	16.58	7.78	3.68	17.86	19.84
						18.33

For word-level EEG features, the impact of pre-train is mostly on BLEU scores while the ROUGE score does not have much variance. Without pre-train, the BLEU- $\{1, 2, 3\}$ respectively drop by 0.64, 1.21, and 1.52. For direct translation on raw waves, the impact is significantly larger. Without self-supervised initialization, the BLEU- $\{1, 2, 3\}$ respectively drop by 3.93 (\downarrow 19.16%), 2.40 (\downarrow 23.57%), and 1.48 (\downarrow 28.68%). Similar observation also appears on ROUGE scores. This phenomenon is rational since raw wave decoding requires the model to pick useful features without any help from eye fixations. The self-supervised initialization could help the model form a preliminary ability to extract time-wise or channel-wise features from raw waves.

5 Limitations

Despite DeWave’s enhancements in EEG-to-Text translation using a discrete codex and raw wave encoding, its accuracy remains far from real-life scenarios compared to traditional language-to-language translations. Also, to keep a fair comparison with EEG-to-Text, this paper uses a teacher-forcing setting in evaluation. This setting eliminates accumulation error and turns the sequence decoding task into a word-level classification task given the ground truth token from the previous step. This setting is relatively easier yet we think it is still valuable as it could suggest feature extraction quality while keeping a fair comparison.

Additionally, the experiments in this paper are restricted to public neural reading data, not fully aligning with the “silent speech” concept of direct thought translation from human brains. Instead, the current ZuCo dataset is collected by giving people reading stimuli. This paper focuses on introducing Wav2Vec formation of feature extraction on raw EEG waves and introduces discrete codex as learnable representations for the EEG-to-Text translation domain. One scientific problem in this domain is a better way of doing the “silent speech” task, which we are exploring as on-going research.

6 Conclusion

This paper presents DeWave, a framework for the recently proposed open-vocabulary EEG-to-Text translation task [48], introducing the concept of discrete codex encoding. This approach brings enhancement in corpus text relevancy metrics, such as BLEU and ROUGE. DeWave also expands the task to decode raw EEG waves without the assistance of eye fixation markers. Despite these advancements, the quality of brain decoding results remains substantially low and remains teacher forcing setting for fair comparison. The translation of thoughts directly from the brain is a valuable yet challenging endeavor that warrants significant continued efforts. In our ongoing work, we are exploring more reasonable settings that remove teacher forcing for both training and testing. We will also include the “neural-feedback” mechanism in this EEG-to-Text research to enhance the scientific value of this domain.

Acknowledgement

This work was supported in part by the Australian Research Council (ARC) under discovery grants DP210101093 and DP220100803, and the GrapheneX-UTS Human-Centric AI Centre sponsored by GrapheneX (2023-2031).

References

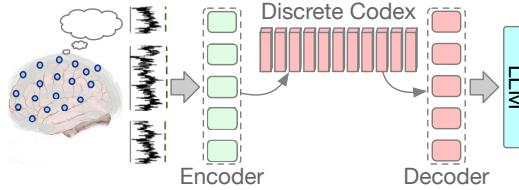
- [1] Gopala K Anumanchipalli, Josh Chartier, and Edward F Chang. Speech synthesis from neural decoding of spoken sentences. *Nature*, 568(7753):493–498, 2019.
- [2] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, 33:12449–12460, 2020.
- [3] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 2020.
- [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [5] Ciaran Cooney, Rafaella Folli, and Damien Coyle. Mel frequency cepstral coefficients enhance imagined speech decoding accuracy from eeg. In *2018 29th Irish Signals and Systems Conference (ISSC)*, pages 1–7. IEEE, 2018.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- [7] Yiqun Duan, Zhen Wang, Yi Li, Jianhang Tang, Yu-Kai Wang, and Chin-Teng Lin. Cross task neural architecture search for eeg signal classifications. *arXiv preprint arXiv:2210.06298*, 2022.
- [8] Yiqun Duan, Zhen Wang, Yi Li, Jianhang Tang, Yu-Kai Wang, and Chin-Teng Lin. Cross task neural architecture search for eeg signal recognition. *Neurocomputing*, 545:126260, 2023.
- [9] Yiqun Duan, Jinzhao Zhou, Zhen Wang, Yu-Cheng Chang, Yu-Kai Wang, and Chin-Teng Lin. Domain-specific denoising diffusion probabilistic models for brain dynamics. *arXiv preprint arXiv:2305.04200*, 2023.
- [10] Jiang et al. Self-supervised contrastive learning for eeg-based sleep staging. In *IJCNN 2021*, pages 1–8. IEEE, 2021.
- [11] Kostas et al. 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.
- [12] Luciano Floridi and Massimo Chiriaci. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30:681–694, 2020.
- [13] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020.
- [14] Christian Herff, Dominic Heger, Adriana De Pesters, Dominic Telaar, Peter Brunner, Gerwin Schalk, and Tanja Schultz. Brain-to-text: decoding spoken phrases from phone representations in the brain. *Frontiers in neuroscience*, 9:217, 2015.
- [15] Nora Hollenstein, Jonathan Rotsztejn, Marius Troendle, Andreas Pedroni, Ce Zhang, and Nicolas Langer. Zuco, a simultaneous eeg and eye-tracking resource for natural sentence reading. *Scientific data*, 5(1):1–13, 2018.
- [16] Nora Hollenstein, Marius Troendle, Ce Zhang, and Nicolas Langer. Zuco 2.0: A dataset of physiological recordings during natural reading and annotation. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asunción Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 138–146. European Language Resources Association, 2020.
- [17] Nora Hollenstein, Marius Troendle, Ce Zhang, and Nicolas Langer. Zuco 2.0: A dataset of physiological recordings during natural reading and annotation. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 138–146, 2020.
- [18] Amir Jahangiri and Francisco Sepulveda. The relative contribution of high-gamma linguistic processing stages of word production, and motor imagery of articulation in class separability of covert speech tasks in eeg data. *Journal of medical systems*, 43(2):1–9, 2019.
- [19] Robert Jenke, Angelika Peer, and Martin Buss. Feature extraction and selection for emotion recognition from eeg. *IEEE Transactions on Affective computing*, 5(3):327–339, 2014.

- [20] Reinmar J Kobler, Jun ichiro Hirayama, Qibin Zhao, and Motoaki Kawanabe. SPD domain-specific batch normalization to crack interpretable unsupervised domain adaptation in EEG. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [21] 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.
- [22] Yi Li, Hualiang Wang, Yiqun Duan, and Xiaomeng Li. Clip surgery for better explainability with enhancement in open-vocabulary tasks. *arXiv preprint arXiv:2304.05653*, 2023.
- [23] Yi Li, Hualiang Wang, Yiqun Duan, Hang Xu, and Xiaomeng Li. Exploring visual interpretability for contrastive language-image pre-training. *arXiv preprint arXiv:2209.07046*, 2022.
- [24] Lun-De Liao, Chi-Yu Chen, I-Jan Wang, Sheng-Fu Chen, Shih-Yu Li, Bo-Wei Chen, Jyh-Yeong Chang, and Chin-Teng Lin. Gaming control using a wearable and wireless eeg-based brain-computer interface device with novel dry foam-based sensors. *Journal of neuroengineering and rehabilitation*, 9(1):1–12, 2012.
- [25] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [26] Joseph G Makin, David A Moses, and Edward F Chang. Machine translation of cortical activity to text with an encoder-decoder framework. *Nature Neuroscience*, 23(4):575–582, 2020.
- [27] Dennis J McFarland, Charles W Anderson, K-R Muller, Alois Schlogl, and Dean J Krusienski. Bci meeting 2005-workshop on bci signal processing: feature extraction and translation. *IEEE transactions on neural systems and rehabilitation engineering*, 14(2):135–138, 2006.
- [28] Anton Nijholt. Bci for games: A ‘state of the art’ survey. In Scott M. Stevens and Shirley J. Saldamarco, editors, *Entertainment Computing - ICEC 2008*, pages 225–228, 2009.
- [29] Julie Onton, Marissa Westerfield, Jeanne Townsend, and Scott Makeig. Imaging human eeg dynamics using independent component analysis. *Neuroscience & biobehavioral reviews*, 30(6):808–822, 2006.
- [30] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [31] Yue-Ting Pan, Jing-Lun Chou, and Chun-Shu Wei. MAtt: A manifold attention network for EEG decoding. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [32] Jerrin Thomas Panachakel, AG Ramakrishnan, and TV Ananthapadmanabha. A novel deep learning architecture for decoding imagined speech from eeg. *arXiv preprint arXiv:2003.09374*, 2020.
- [33] Jerrin Thomas Panachakel and Angarai Ganesan Ramakrishnan. Decoding covert speech from eeg-a comprehensive review. *Frontiers in Neuroscience*, 15:392, 2021.
- [34] 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.
- [35] 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. *arXiv preprint arXiv:2103.00020*, 2021.
- [36] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [37] Jaime F Delgado Saa and Müjdat Çetin. Discriminative methods for classification of asynchronous imaginary motor tasks from eeg data. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 21(5):716–724, 2013.
- [38] David Sabbagh, Pierre Ablin, Gael Varoquaux, Alexandre Gramfort, and Denis A. Engemann. Manifold regression to predict from meg/eeg brain signals without source modeling. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [39] Pramit Saha, Muhammad Abdul-Mageed, and Sidney Fels. Speak your mind! towards imagined speech recognition with hierarchical deep learning. *arXiv preprint arXiv:1904.05746*, 2019.
- [40] Simanto Saha and Mathias Baumert. Intra-and inter-subject variability in eeg-based sensorimotor brain computer interface: a review. *Frontiers in computational neuroscience*, 13:87, 2020.

- [41] Jingyuan Sun, Shaonan Wang, Jiajun Zhang, and Chengqing Zong. Towards sentence-level brain decoding with distributed representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7047–7054, 2019.
- [42] Pengfei Sun, Gopala K Anumanchipalli, and Edward F Chang. Brain2char: a deep architecture for decoding text from brain recordings. *Journal of Neural Engineering*, 17(6):066015, 2020.
- [43] Kazuo Tanaka, Kazuyuki Matsunaga, and Hua O Wang. Electroencephalogram-based control of an electric wheelchair. *IEEE transactions on robotics*, 21(4):762–766, 2005.
- [44] Luca Tonin, Tom Carlson, Robert Leeb, and José del R. Millán. Brain-controlled telepresence robot by motor-disabled people. In *2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 4227–4230, 2011.
- [45] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [46] Neeraj Wagh, Jionghao Wei, Samarth Rawal, Brent M. Berry, and Yogatheesan Varatharajah. Evaluating latent space robustness and uncertainty of EEG-ML models under realistic distribution shifts. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [47] Jiang Wang and Mei Wang. Review of the emotional feature extraction and classification using eeg signals. *Cognitive Robotics*, 1:29–40, 2021.
- [48] 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.
- [49] 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, 2021.
- [50] Haoning Xi, Didier Aussel, Wei Liu, S Travis Waller, and David Rey. Single-leader multi-follower games for the regulation of two-sided mobility-as-a-service markets. *European Journal of Operational Research*, 2022.
- [51] Haoning Xi, Liu He, Yi Zhang, and Zhen Wang. Differentiable road pricing for environment-oriented electric vehicle and gasoline vehicle users in the bi-objective transportation network. *Transportation Letters*, 14(6):660–674, 2022.
- [52] Xiang Zhang, Lina Yao, Quan Z Sheng, Salil S Kanhere, Tao Gu, and Dalin Zhang. Converting your thoughts to texts: Enabling brain typing via deep feature learning of eeg signals. In *2018 IEEE international conference on pervasive computing and communications (PerCom)*, pages 1–10. IEEE, 2018.
- [53] Shunan Zhao and Frank Rudzicz. Classifying phonological categories in imagined and articulated speech. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 992–996. IEEE, 2015.
- [54] Jinzhao Zhou, Yiqun Duan, Yu-Cheng Chang, Yu-Kai Wang, and Chin-Teng Lin. Belt: Bootstrapping electroencephalography-to-language decoding and zero-shot sentiment classification by natural language supervision. *arXiv preprint arXiv:2309.12056*, 2023.
- [55] Jinzhao Zhou, Yiqun Duan, Zhihong Chen, Yu-Cheng Chang, and Chin-Teng Lin. Generalizing multimodal variational methods to sets, 2022.
- [56] Jinzhao Zhou, Yiqun Duan, Yingying Zou, Yu-Cheng Chang, Yu-Kai Wang, and Chin-Teng Lin. Speech2eeg: Leveraging pretrained speech model for eeg signal recognition. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 2023.

Supplementary Material for DeWave: Discrete Encoding of EEG Waves for EEG to Text Translation

In this material, we will give more technical details as well as additional experiments to support the main paper. The overview of the proposed framework, DeWave, is illustrated in Figure 6.



Ground Truth	Bush attended <u>the University of Texas at Austin</u> , where he graduated Phi Beta Kappa with a Bachelor's degree in <u>Latin American Studies in 1973</u> , taking only two and a half years to complete his work, and obtaining generally excellent grades .
Predict	was <u>the University of California at Austin</u> in where he studied in Beta Kappa in a degree of degree in <u>history American Studies in 1975</u> . and a one classes a half years to complete the degree. and was a excellent grades .

Figure 6: Overall illustration of translating EEG waves into text through quantised encoding.

A Dataset

ZuCo stands for Zurich Cognitive Language Processing Corpus (ZuCo), a dataset that includes both raw and preprocessed eye-tracking and electroencephalography (EEG) data. The data is collected by having human subjects read given text corpora while simultaneously recording both their eye-tracking signals and EEG waves. The recording is done using the Biosemi-128 system, which, after denoising, provides 105 out of 128 channels for downstream tasks. The dataset comprises two versions: ZuCo 1.0, collected from 12 subjects, and ZuCo 2.0, collected from 18 subjects [15, 17].

The text corpora within the ZuCo dataset are sourced from a diverse set of textual genres, including 1) Wikipedia articles, 2) movie reviews, and 3) the BNC (British National Corpus). This diversity ensures a wide variety of syntactic structures and word frequencies. The dataset records data during two tasks: Normal Reading (NR) and Task-Specific Reading (TSR). In our experiments, DeWave utilizes both ZuCo 1.0 [15] and 2.0 [17]. The EEG features are captured using a 128-channel system with a sampling rate of 500Hz, filtered through a frequency band ranging from 0.1Hz to 100Hz. After noise canceling, only 105 channels are deemed suitable for translation [15].

For word-level EEG feature translation, eye-fixation data associated with each word during reading is available in the ZuCo dataset. Following the approach similar to [48], we extract segments of the EEG wave according to eye fixations. Words fixated upon multiple times have their EEG fragments concatenated for processing. To process these word-level EEG features, we compute statistical results across four frequency band filters: the Theta band (5-7Hz), the Alpha band (8-13Hz), the Beta band (12-30Hz), and the Gamma band (30Hz and above) [27]. Consequently, the feature size for each word totals $105 \times 4 \times 2 = 840$. For raw EEG waves, the signals are normalized to a range between 0 and 1 for decoding.

The dataset is split into training (80%), development (10%), and testing (10%) sets, comprising 10,874, 1,387, and 1,387 unique sentences, respectively, with no overlap. We further conducted a statistical analysis on the sentences extracted from the dataset, details of which are reported below.

Table 4: Statistical analysis of sentences from the ZuCo dataset.

Feature	ZuCo 1.0 Natural Reading	ZuCo 2.0 Natural Reading
Sentences	300	390
Sent. length	21.3 ± 10.6	19.6 ± 8.8
Total words	6386	6828
Word length	6.7	4.9

B Implementation Details

We release our implementation code through GitHub to contribute to this area. Currently the basic code are available through an anonymous link⁶ For word-level EEG features, we use the 56 tokens each with an 840 embedding size. The codex encoder for word-level features is a 6-layer transformer encoder with head number 8, hidden embedding 512. For raw EEG waves, we clip or pad the EEG waves up to sample point 5500 with a constant value of zero, which scales up to 11 seconds according to the sampling rate of 500 Hz. The codex encoder for raw EEG wave features is illustrated in Section 3.3, where a 6-layer CNN encoder slides through the whole wave and gets the embedding sequence. A transformer layer with head number 8 and a 1×1 convolutional layer are combined to fuse multiple EEG channels into one embedding with size 512.

The codex encoder shares the same structure with word-level features. DeWave uses a codex with size 2048 where each codex latent is an embedding with size 512. The ablation study gives a discussion about the codex size. All models are trained on Nvidia V100 and A100 GPUs. For the self-supervised decoding for raw waves, we use a learning rate of 5e-4 and a VQ coefficient of 0.25 for training 35 epochs. For training the codex (stage 1), DeWave uses a learning rate of 5e-4 for 35 epochs. For finetuning the translation (stage 2), DeWave uses a learning rate of 5e-6 for 30 epochs. We use the SGD as the optimizer for training all the models.

C Training Paradigm

DeWave is trained through a multi-stage process, where the training process is illustrated in Appendix algorithm 1. Before the two-stage training, if the input of the model is the raw waves, we initialize the wave2vec model with a self-supervised pre-training described in section 3.3. The self-supervised training is realized by encoding raw waves into discrete codex and reconstructing the discrete codex into original raw waves. The training process for self-supervised initialization utilizes the SGD algorithm with a learning rate of 0.0005 for 30 epochs with 0.1 times the learning rate decrease at epoch 20. In the first stage, we do not involve the language model in weight updates. The target of the first stage is to train a proper encoder projection θ_{codx} and a discrete codex representation \mathcal{C} for the language model. Intuitively, if the learning of the codex is successful, the translator could receive a better representation that is closer to the original representation, word2vec embedding. The training for the first stage is optimized by the SGD optimizer with a learning rate of 0.0005 for 35 epochs. However, the language model is trained on word tokens, which may not be perfectly suitable for brain tokens. In the second stage, the gradient of all weights, including language model θ_{BART} is opened to fine-tune the whole system. The training for the second stage is optimized by the SGD optimizer with a learning rate of 1e - 6 for 35 epochs.

D Codex Visualization

Since the logic is to learn a discrete codex from brain dynamics, it naturally arises whether the codex value distribution between raw waves and frequency features has differences. In that case, we conduct additional experiments to visualize the learned codex book with T-SNE methods and report the results in Fig 7. Ideally, the purpose of the discrete codex is to make the language model have a better understanding of brain encoding. In that case, the learned codex regardless of whether it is for frequency features or raw waves, should be approximately the same as the word2vec embedding. In other words, the distribution of the codex should be similar. Fig. 7 supports our expectation,

⁶<https://github.com/duanyiqun/DeWave>

Algorithm 1 Training procedure

Input: EEG \mathcal{E} , Vocabulary \mathcal{V} , Marker \mathcal{F} , Target \mathcal{W}
Parameter: Codex $\mathcal{C} = \{\mathbf{c}\}$, θ_{codex} , θ_{BART} , Θ_{wave}

```
1: if decode raw waves then
2:    $\arg \min_{\mathcal{C}, \theta_{codex}, \Theta_{wave}} L_{wave}$ 
3:   Vectorize  $\mathcal{X} = \Theta(\mathcal{E})$ 
4: else
5:   Vectorize  $\mathcal{X} = \Theta(\mathcal{E}, \mathcal{F})$ 
6: end if
7: Stage-1: Train codex
8: while Iteration steps do
9:    $\arg \min_{\mathcal{C}, \theta_{codex}} L(\mathcal{X})$ 
10: end while
11: Stage-2: Finetune Language Model
12: while Iteration steps do
13:    $\arg \min_{\mathcal{C}, \theta_{codex}, \theta_{BART}} L(\mathcal{X})$ 
14: end while
15: return  $\mathcal{C}, \theta_{codex}, \theta_{BART}, \Theta_{wave}$ 
```

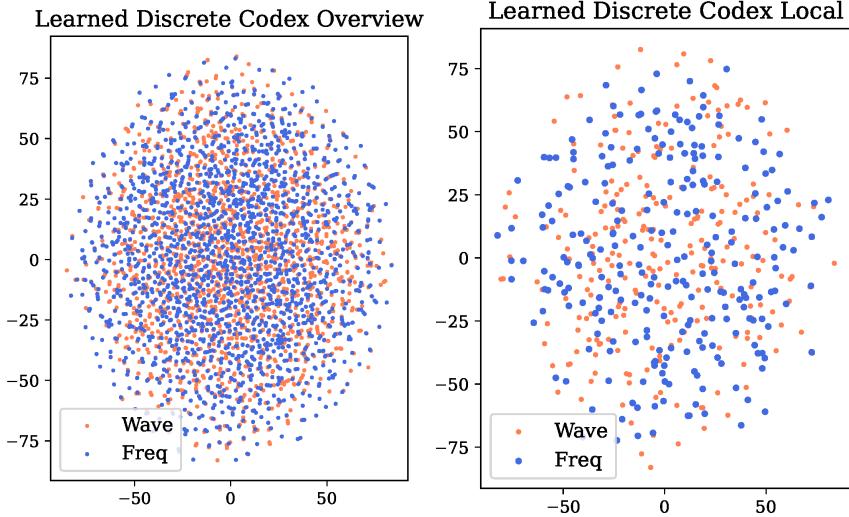


Figure 7: Visualization of codex value distribution, where the left is the global distribution, and the right is the local distribution.

where the codex learned from frequency and raw waves have very similar distributions. However, the frequency codex has more coverage with corner cases and boundaries. We think this is rational since frequency features are naturally easier to distinguish as it has introduced manually selected features. The performance gap between raw waves and frequency features supports this point as well. Still, the similarity of the distribution illustrates the rationality of the learned codex.

E Motivation and Preliminary Tests with LLMs

We provide additional insights into our experiments with larger language models. While our primary experiments utilized BART to ensure consistency in the decoder scale with prior works, it was paramount for us to ascertain that the observed improvements emanated from discrete coding and not just a more sophisticated decoder.

The potential of bridging brain activities with larger language models (LLMs) and advancing towards AGI is a significant research avenue. Recognizing this, we recently undertook an ablation study, where we replaced the BART decoder with OPT and Llama V1. Contrary to our expectations, the performance enhancement was modest. Given the vast implications of this area, we previously refrained from including these findings in the main manuscript for reasons of prudence.

Limited by our computational resources, we employed the PyTorch FSDP mode to fine-tune the OPT-1.3B and Llama-1 7B models with half-precision across three epochs. Taking cues from Mini-GPT4’s method for handling visual tokens, the tokenized EEG waves were prompted into LLMs. The performance metrics for our experiments with BART, OPT 1.3B, and Llama-1 7B decoders are tabulated below:

Table 5: Performance metrics for different decoders on the ZuCo dataset.

Source	Decoder	BLEU-1	BLEU-3	ROUGE-R	ROUGE-P	ROUGE-F
Word-level features	DeWave	41.35	13.92	28.82	33.71	30.69
Word-level features	DeWave + OPT 1.3B	41.97	14.06	28.98	33.82	30.86
Word-level features	DeWave + Llama-1 7B	42.84	15.03	29.42	35.43	32.05
Raw Waves	DeWave	20.51	5.16	21.18	29.42	24.27
Raw Waves	DeWave + OPT 1.3B	21.31	5.84	22.09	29.94	25.42
Raw Waves	DeWave + Llama-1 7B	22.05	6.03	22.45	30.01	26.08

From the table, it’s evident that while the LLMs offer some enhancement, the gains are not as pronounced as one might expect. This underscores the complexity of the problem and the challenges of bridging brain activities with LLMs. This simple experiment provides a deeper dive into our experiments with larger language models. We believe these findings offer additional perspectives and pave the way for more nuanced research in this domain.

F Generated Samples

In this section, we visualize the generated decoding text results on brain waves and compare them with the ground truth in Table 6 and Table 7. It suggests that the results are even better on long and simple sentences. For example, even for the ground truth with long and logic as below, the prediction could still match key information throughout the whole sentence.

Ground Truth:

Bush attended the University of Texas at Austin, where he graduated Phi Beta Kappa with a Bachelor’s degree in Latin American Studies in 1973, taking only two and a half years to complete his work, and obtaining generally excellent grades.

Model Output:

was the University of California at Austin in where he studied in Beta Kappa in a degree of degree in history American Studies in 1975. and a one classes a half years to complete the degree. and was a excellent grades.

People is feasible to guess the meaning of a human based on the translation from brain waves. In the example above, the model recognizes through waves that the University of xxx at Austin. Although it is a factual mistake that the University of California is not at Austin, it still suggests that the model could approximately capture the semantic meaning through non-invasive brain waves. A similar observation applies that the model recognizes that it is the xx American Studies in xx years however the model predicts history American Studies in 1975 rather than the ground truth is Latin American Studies in 1973. Surprisingly, even the years have correlations at this stage.

Table 6: Translation comparison between the ground truth and the prediction on brain waves with eye fixation on task v2.0 dataset.

(1)	Ground Truth: The book was awarded the 1957 Pulitzer Prize for Biography ... Prediction: first is published the Pulitzer Pulitzer Prize for Literatureography ...
(2)	Ground Truth: Kennedy's other decorations of the Second World War include the Purple Heart, Asiatic-Pacific Campaign Medal, and the World War II Victory Medal. Prediction: eth was son son were the day World War were a famous Heart and thepenatic StarAmerican,,, and the American War II Victory Medal.
(3)	Ground Truth: In 1958, Kennedy published the first edition of his book A Nation of Immigrants, closely following his involvement in the Displaced Persons Act and the 1957 bill to bring families together. Prediction: the, the was his novel of of his autobiography, Life of Millionsigrants, which followed the experiences in the Vietnamrael Persons Movement of the Civil assassination of abolish it together.
(4)	Ground Truth: After World War II, Kennedy entered politics (partly to fill the void of his popular brother, Joseph P. Kennedy, Jr., on whom his family had pinned many of their hopes but who was killed in the war) ... Prediction: the War II, the was the asasly as avoid a void left a father father, John Kennedy. Kennedy), who.) who the he father had been the hopes the hopes). who had assassinated in the war. ...
(5)	Ground Truth: In 1946, Representative James Michael Curley vacated his seat in an overwhelmingly Democratic district to become mayor of Boston and Kennedy ran for that seat, beating his Republican opponent by a large margin. Model Output: the, the John W Smithley was the seat in the unsuccessful Republican Congress of become a of New. become's for president office in which incumbent opponent opponent, a landslide margin.
(6)	Ground Truth: He was reelected twice, but had a mixed voting record, often diverging from President Harry S. Truman and the rest of the Democratic Party. Model Output: was a- to in in lost to less record record. and votingting from the Obama Truman. Truman's his Republican of the Republican Party.
(7)	Ground Truth: He was reelected twice, but had a mixed voting record, often diverging from President Harry S. Truman and the rest of the Democratic Party. Model Output: was a- to in in lost to less record record. and votingting from the Obama Truman. Truman's his Republican of the Republican Party.
(7)	Ground Truth: However, the U.S. Navy accepted him in September of that year. Model Output: it film.S. government has the as the. that year.
(8)	Ground Truth: In the spring of 1941, Kennedy volunteered for the U.S. Army, but was rejected, mainly because of his troublesome back. Model Output: the meantime of 2016, the was to the first.S. Army. and was discharged for and because he his age temper.
(9)	Ground Truth: When Bush was seventeen, he went to Leon, Mexico, as part of his school's student exchange program. Model Output: the was president, he was to aidas Nebraska, to a of a father's " exchange program.
(10)	Ground Truth: In November 1977 he was sent to the Venezuelan capital of Caracas, in South America, to open a new operation for the bank. Model Output: the,, was born to the United prison, Manacas to where a America, to work a restaurant bank. the government.
(11)	Ground Truth: In 1923 he was awarded the inaugural Bôcher Memorial Prize by the American Mathematical Society. Model Output: the, was born the Nobel PulitzerAFTAn Prize Medal for the French Academyical Society.
(12)	Ground Truth: The mathematician Garrett Birkhoff (1911-1996) was his son. Model Output: tfirst and Wkoff was1802-19) was a name.

Table 7: Translation comparison between the ground truth and the prediction on brain waves with eye fixation on task v2.0 dataset.

(13)	Ground Truth: Jeb Bush was born in Midland, Texas, where his father was running an oil drilling company. Model Output: uan Bush was a in 18way, Texas, in he father was an insurance refinery company.
(14)	Ground Truth: He was noted for his lyrical playing, and performed with John Coltrane, Dexter Gordon, Hampton Hawes, Jackie McLean, and Ike and Tina Turner, among others. Model Output: was a for his "ical, style and his a a Legendtrane in who Gordon and and Fes and and GleGovern and and others Turner Tina Turner. among others.
(15)	Ground Truth: He later became an educator, teaching music theory at the University of the District of Columbia; he was also director of the District of Columbia Music Center jazz workshop band. Model Output: was added a actor and and at to and the University of California Arts of Columbia. he was also a of the University's Columbia's Festival. program..
(16)	Ground Truth: John Ellis "Jeb" Bush (born February 11, 1953), a Republican, is the forty-third and current Governor of Florida. Prediction: nie,Johnnock" Ellis (19 18 17, 18) a former, was the author-six president final governor of Texas.
(17)	Ground Truth: He is a prominent member of the Bush family, the younger brother of President George W. Bush and the second son of former President George H. W. Bush and Barbara Bush. Prediction: was a former member of the American family and and first brother of President Bush Bush. Bush. the father son of President President Richard H. W. Bush. his Bush.
(18)	Ground Truth: After earning his degree, Bush went to work in an entry level position in the international division of Texas Commerce Bank, which was run by Ben Love. ... Prediction: the his degree, he was on work for the office- position in the Department banking of the A.. where was later by his Carsonll ...
(19)	Ground Truth: Following the 1980 presidential election, Bush and his family moved to Miami-Dade County, Florida. Model Output: the deaths election, the was his wife moved to California,Dade County, Florida,
(20)	Ground Truth: He took a job in real estate with Armando Codina, a 32-year-old Cuban immigrant and self-made American millionaire. Model Output: was a liking as the estate in aando Iino in a former-year-old from immigrant from former-made millionaire businessman.
(21)	Ground Truth: [4] Situated in Liberty City, Dade County, the school is located just outside of greater Miami, in an area plagued by poverty. Model Output: ..]] Theuations in the,, Missouri. County, North city is a in outside of the New. Florida the area known by crime and
(22)	Ground Truth: The co-founder, working alongside Bush as a partner, was T. Williard Fair, a well-known local black activist and head of the Greater Miami Urban League. Model Output: first-founder of John with his, a consultant, was aoniJard,banks who former-known film politiciansmith and activist of the Black Chicago NAACP League.
(23)	Ground Truth: Governor Buddy MacKay (55% to 45%) to become governor, after courting moderate voters and Hispanics. Model Output: or of RoKay ofleft) of 55%) of the governor of and theting the opposition in winning.
(24)	Ground Truth: At the urging of his wife, Columba, a devout Mexican Catholic, the Protestant Bush became a Roman Catholic. Model Output: the same of his wife, hea, he young Catholic Catholic, he actor pastorman a Catholic Catholic in
(25)	Ground Truth: Bush attended the University of Texas at Austin, where he graduated Phi Beta Kappa with a Bachelor's degree in Latin American Studies in 1973, taking only two and a half years to complete his work, and obtaining generally excellent grades. Model Output: was the University of California at Austin in where he studied in Beta Kappa in a degree of degree in history American Studies in 1975. and a one classes a half years to complete the degree. and was a excellent grades.

G Subject Wise Evaluation

Section 4.3 introduces subject-wise metric evaluation on word-level EEG features by removing the subject to be tested from the training data, and then training the model from scratch for testing. The results are shown in Fig 4 in the main paper, where different subjects share the same reading article. However, in supplementary details, we conduct a more detailed subject-wise evaluation in that we respectively train the model on every single subject on task v2.0 and test every subject to report the cross-subject performance. The single subject denotes the model only trained on a single subject on the task v2.0 dataset. For each subject, however, due to limited data, we add all data from task 1.0 as an assistance base. The results for each subject are reported below.

Here we selected subject YAC (Table 8), YAK (Table 9), YDG (Table 10), YFS (Table 11), YSL (Table 12), and YMD (Table 13) to report the performance. We visualize the metrics by clustering the same metrics value of different subjects in one radar chart. It is observed that the model performance might not be optimized if we train and test on the same subject. For example, if we train the model with task2.0 data from subject YFS and test on all subjects, the YFS subject only reaches BLUE-{1 – 4} 43.32, 26.30, 14.68, and 7.62, which is lower than YDS, YAG, YRP, .. etc. which reach 43.79, 26.46, 14.94, and 8.03. This suggests the cross-subject robustness of the proposed DeWave model. Also, since we use the same visualization scale for each radar chart, the area of the chart suggests the performance level. It is observed that if we change the training data from subject to subject, the average performance of every subject is affected by a similar trend. We think this phenomenon is caused by the different signal-to-noise ratios. Some subjects might naturally have less noise interference which makes it easier for the model to learn meaningful features during the training process.

H Supplementary Conclusion

In this supplementary material, we give implementation details, training schema, and most importantly, more generated results and subject-wise evaluation of the proposed DeWave model. The generated results suggest a surprisingly good correlation between the model output on brain waves and the ground truth, even in long sentences with logic. Although there are factual and cognitive mistakes in the translation, it is still feasible to guess the meaning of a human based on the translation from brain waves. The subject-wise evaluation suggests the DeWave model is stable across different human subjects. Please refer to the tables attached below.

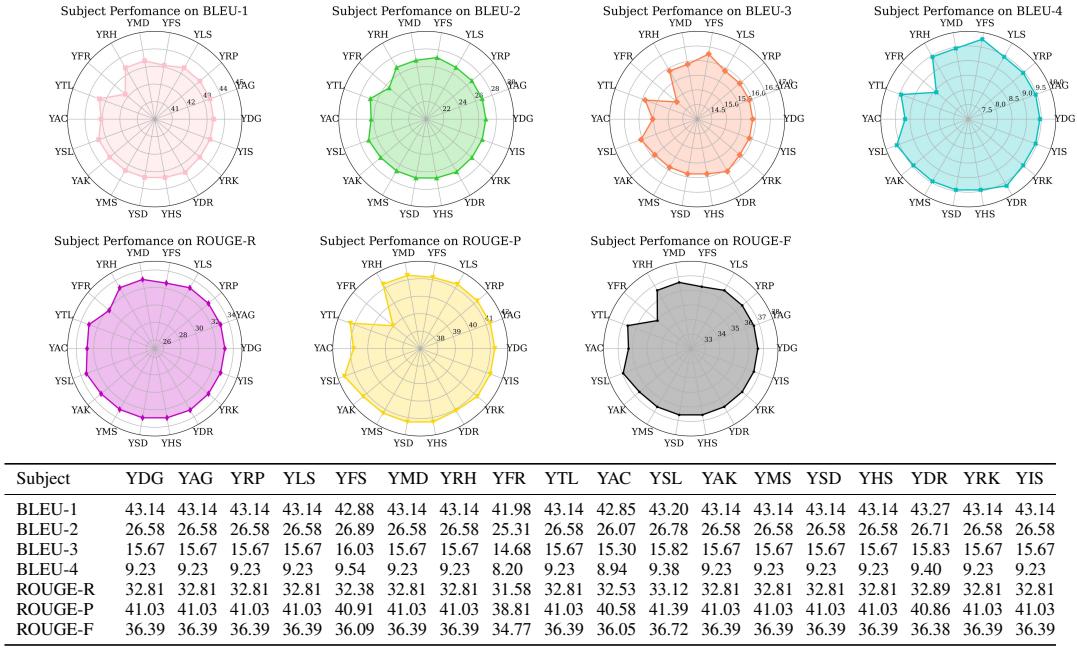


Table 8: Subject-wise evaluation results on a model trained with subject **YAC**, where the radar chart suggests the performance variance on different subjects on each metric.

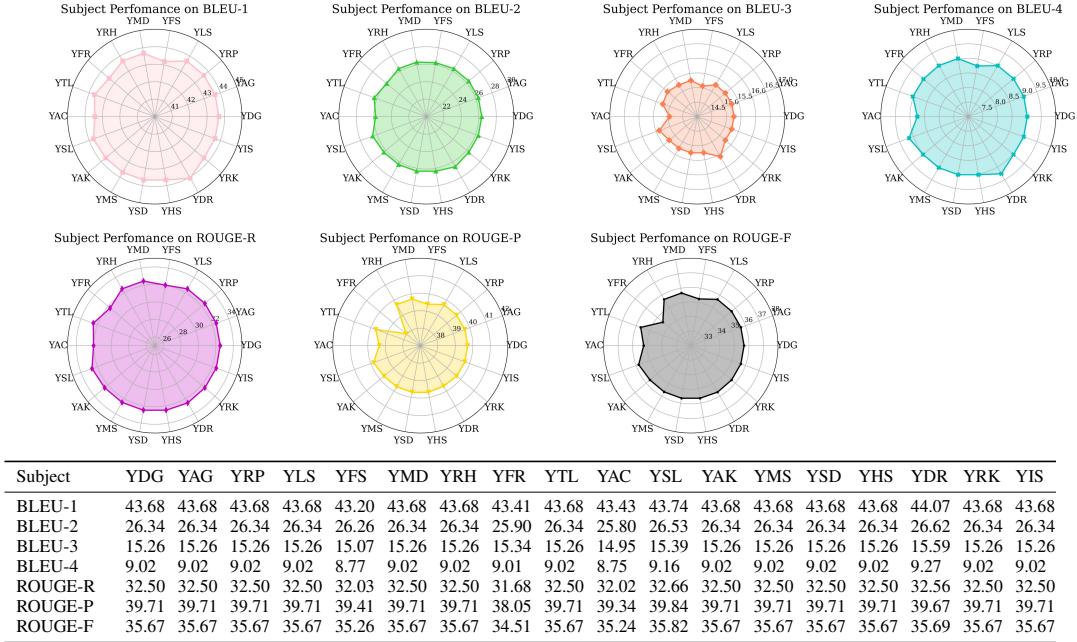


Table 9: Subject-wise evaluation results on a model trained with subject **YAK**, where the radar chart suggests the performance variance on different subjects on each metric.

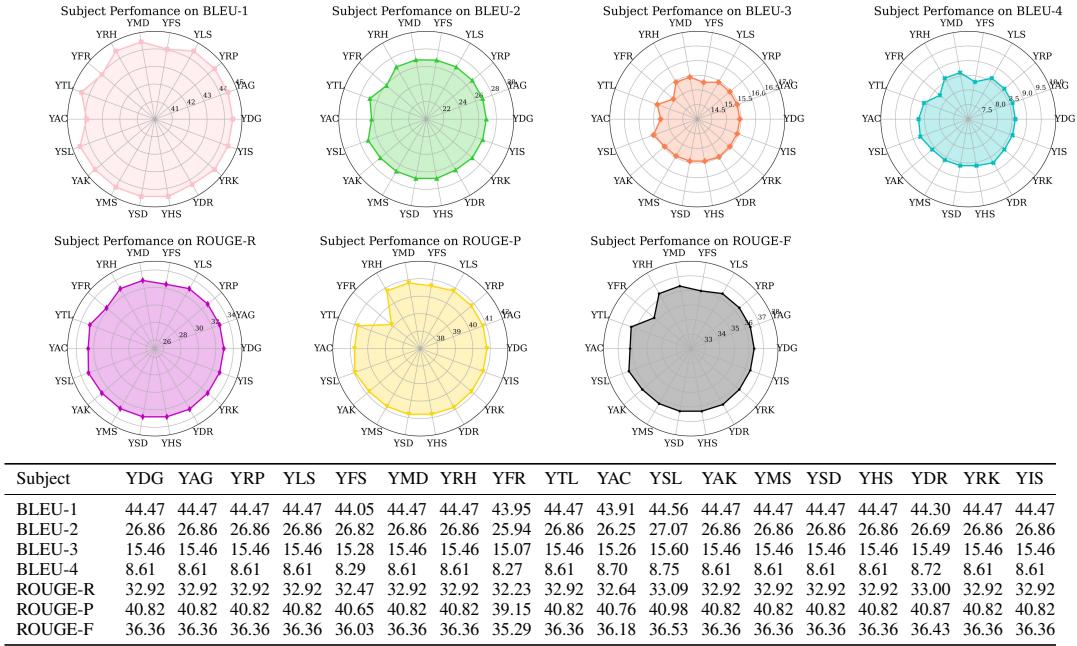


Table 10: Subject-wise evaluation results on a model trained with subject **YDG**, where the radar chart suggests the performance variance on different subjects on each metric.

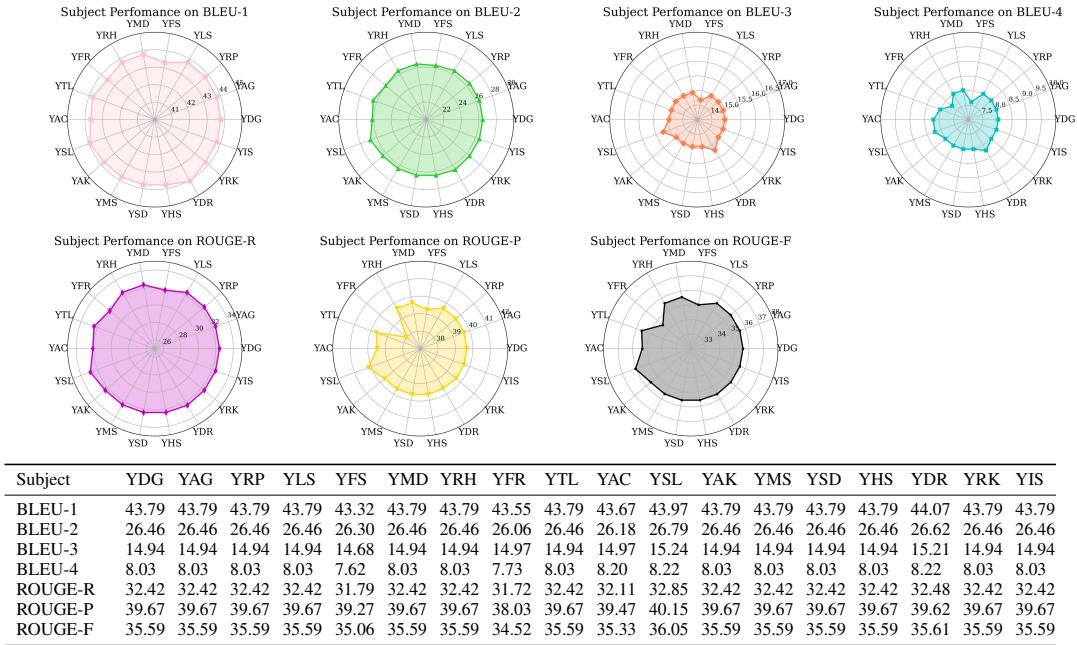


Table 11: Subject-wise evaluation results on a model trained with subject **YFS**, where the radar chart suggests the performance variance on different subjects on each metric.

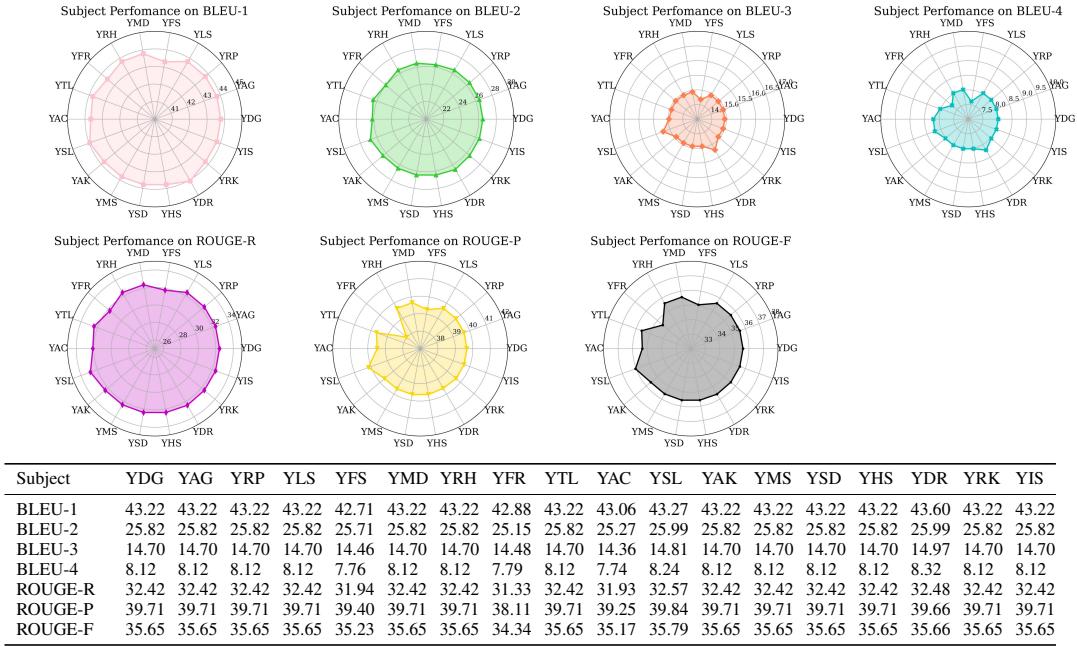


Table 12: Subject-wise evaluation results on a model trained with subject **YSL**, where the radar chart suggests the performance variance on different subjects on each metric.

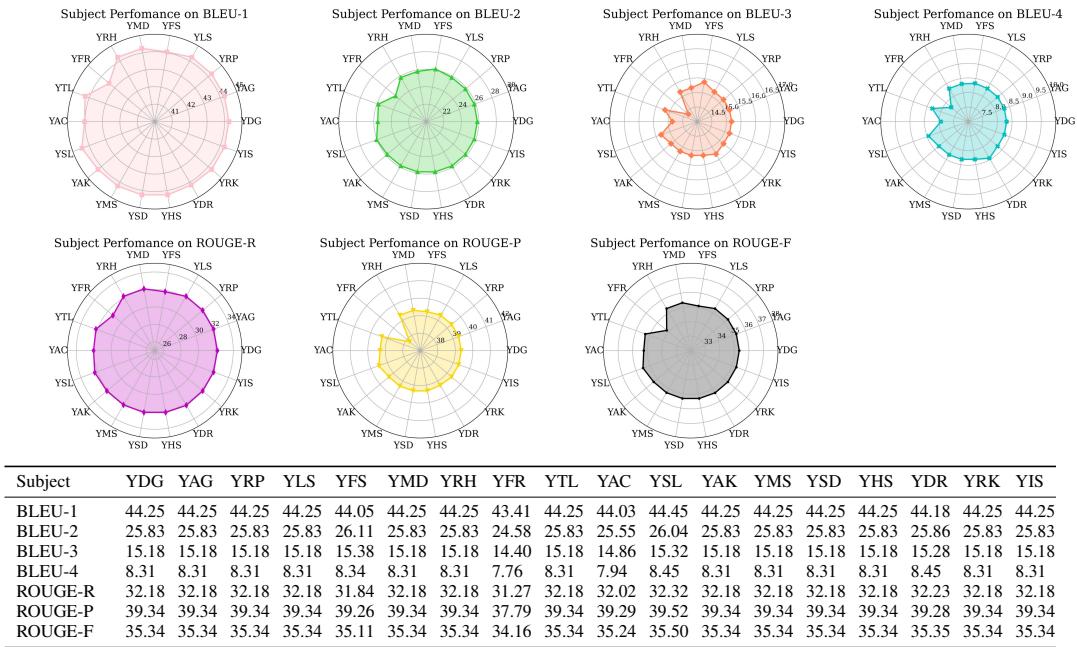


Table 13: Subject-wise evaluation results on a model trained with subject **YMD**, where the radar chart suggests the performance variance on different subjects on each metric.