# Brain-Driven Representation Learning Based on Diffusion Model

Soowon Kim
*Dept. of Artificial Intelligence*
*Korea University*
Seoul, Republic of Korea
soowon_kim@korea.ac.kr

Seo-Hyun Lee
*Dept. of Brain and Cognitive Engineering*
*Korea University*
Seoul, Republic of Korea
seohyunlee@korea.ac.kr

Young-Eun Lee
*Dept. of Brain and Cognitive Engineering*
*Korea University*
Seoul, Republic of Korea
ye_lee@korea.ac.kr

Ji-Won Lee
*Dept. of Artificial Intelligence*
*Korea University*
Seoul, Republic of Korea
jiwon_lee@korea.ac.kr

Ji-Ha Park
*Dept. of Artificial Intelligence*
*Korea University*
Seoul, Republic of Korea
jiha_park@korea.ac.kr

Seong-Whan Lee
*Dept. of Artificial Intelligence*
*Korea University*
Seoul, Republic of Korea
sw.lee@korea.ac.kr

*Abstract*—Interpreting EEG signals linked to spoken language presents a complex challenge, given the data's intricate temporal and spatial attributes, as well as the various noise factors. Denoising diffusion probabilistic models (DDPMs), which have recently gained prominence in diverse areas for their capabilities in representation learning, are explored in our research as a means to address this issue. Using DDPMs in conjunction with a conditional autoencoder, our new approach considerably outperforms traditional machine learning algorithms and established baseline models in accuracy. Our results highlight the potential of DDPMs as a sophisticated computational method for the analysis of speech-related EEG signals. This could lead to significant advances in brain-computer interfaces tailored for spoken communication.

*Keywords–brain-computer interface, electroencephalogram, imagined speech, diffusion model;*

## I. INTRODUCTION

Speech serves as an essential means of human communication, allowing us to express intricate thoughts and ideas through audible patterns. Speech capacity is deeply embedded in our social and cultural fabric, facilitating everything from relationship building to information sharing. Despite its importance, some people, such as those suffering from locked-in syndrome, are unable to engage in verbal communication due to physical limitations [1]. Therefore, innovative approaches to restore or replace speech capabilities are a vital research frontier. In line with this, our work focuses on the interpretation of brain signals as a means for facilitating non-vocal communication.

Electroencephalography (EEG) offers a non-invasive avenue for capturing the brain's electrical activities. Acquired through scalp-placed electrodes, EEG signals are instrumental in various applications, ranging from neuroscience to clinical diagnostics [2]. The decoding of these EEG signals into useful data, such as speech-related activities or focus levels, is of growing interest.

Deciphering EEG data related to spoken language is notably intricate. The task involves interpreting complex and highly variable neural activities related to the articulation and perception of speech. Additionally, these EEG signals often contain noise and artifacts, further complicating accurate decoding. In light of these challenges, ongoing research aims to establish robust and effective methods for EEG signal interpretation, which have broad applications including speech restoration and human-machine interactions.

Denoising diffusion probabilistic models (DDPMs) have emerged as a potent tool for identifying nuanced patterns within complicated, high-dimensional datasets. Through a process of adding Gaussian noise over a series of steps, DDPMs corrupt an original signal and then attempt to reconstruct it. These models have been particularly successful in dealing with time series data, including audio and video streams.

Decoding EEG signals using deep learning approaches is a challenging problem due to various factors, including the scarcity of data, a poor signal-to-noise ratio, and high inter- and intra-individual variability [3]. Despite these challenges, several studies have explored different EEG decoding techniques for various applications, including speech decoding [4].

On the basis of existing research, various methods for the decoding of EEG signals have been explored. Schirrmeister et al. utilized DeepConvNets [5] to achieve end-to-end learning in human EEG signals, using machine learning techniques

such as batch normalization and exponential linear units. They achieved performance on par with traditional filter bank common spatial pattern algorithms. Lawhern et al. [6], [7] introduced EEGNet, a specialized CNN architecture for EEG classification, using depthwise and separable convolutions to better capture specific EEG features. This model has been successfully tested on multiple BCI paradigms.

Furthermore, Lee et al. [8], [9] conducted an in-depth study of the nuances that affect decoding performance in two key BCI paradigms: imagined speech and visual imagery. The study used EEG signals filtered across multiple frequency ranges and identified relevant cortical regions [10], resulting in high precision and multiclass scalability for both paradigms.

In parallel, diffusion-based approaches for time series data have gained substantial traction. One such method, proposed by Alcaraz et al., uses a structured state-space model with an integrated diffusion process for time-series data imputation and forecasting, showing superior performance to existing methods. Jeong et al. offered a novel application of diffusion models to improve synthetic speech quality in Text-to-Speech (TTS) systems, which also demonstrated effectiveness over current methods.

Our study builds on this background to introduce a new strategy for interpreting EEG signals linked to spoken language using DDPM and a conditional autoencoder (CAE). The CAE facilitates the retention of valuable features that may otherwise be compromised during the DDPM's forward process. Additionally, we incorporate a jointly trained classifier to enhance decoding performance. To the best of our knowledge, this is the first effort to apply diffusion models to interpret speech-related EEG signals.

Our study extends this body of work by introducing a novel approach that combines DDPMs and a conditional autoencoder to decode EEG signals related to spoken language. This method aims to capture the intricate neural patterns and relationships inherent in speech processes, and, in doing so, advances the field of EEG decoding with potential applications in speech rehabilitation and brain-computer interfaces.

## II. MATERIALS AND METHODS

### A. Denoising Diffusion Models

DDPMs are a type of machine learning model that can learn complex probability distributions over data. The "forward process" in DDPMs is determined by a fixed Markov chain that progressively adds Gaussian noise to the data. The forward process begins with a probability distribution, denoted $q(\mathbf{x}_0)$, which represents the uncorrupted original data. This distribution is then iteratively transformed using a sequence of Markov diffusion kernels, $q(\mathbf{x}_t \mathbf{x}_{t-1})$, which are Gaussian with a fixed variance schedule $\{\beta_t\}_{t=1}^T$. This process can be expressed as follows:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad (1)$$
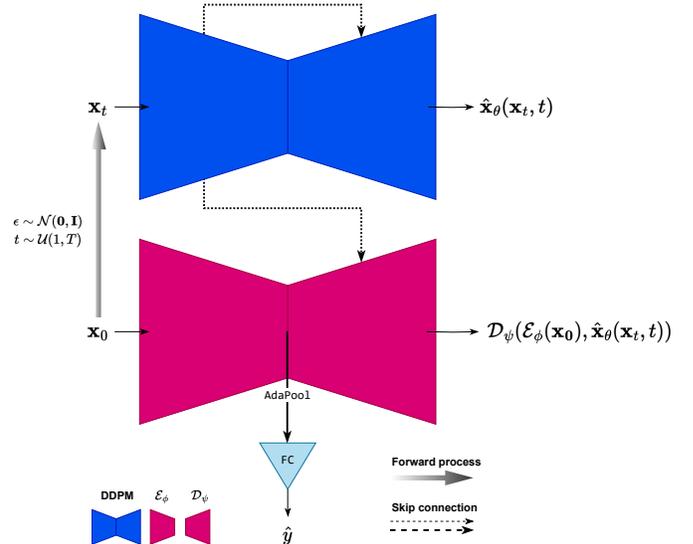


Fig. 1. A flowchart of Diff-E for EEG signal decoding. Initially, DDPM processes noisy data to approximate the original signal, then a CAE refines this output by correcting the discrepancies. Subsequently, the classifier utilizes the encoder's output for downstream classification tasks, enhancing overall performance.

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}). \quad (2)$$

The original data can be corrupted by a diffusion process with Gaussian noise at any stage, $t$, where $\alpha_t$ is represented as $1 - \beta_t$ and $\bar{\alpha}_t$ is the product of all $\alpha_s$ from $s = 1$ to $t$.

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}). \quad (3)$$

Ho et al. [8] proposed a technique to train a model that takes a noisy sample $\mathbf{x}_t$ and predicts the noise it contains by training a network $\epsilon\theta(\mathbf{x}_t, t)$. On the contrary, our research trains Diff-E to forecast the original unchanged signal, $\mathbf{x}_0$, rather than predicting the injected noise.

$$\mathcal{L}_{\text{DDPM}}(\theta) = ||\mathbf{x}_0 - \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)||. \quad (4)$$

We randomly select a timestep $t$ from a uniform distribution, $\mathcal{U}(1, T)$, and use $\theta$ as the parameters of the DDPM. The objective of the model is to denoise the noisy input and generate an output that is close to the original signal. We have employed a time-conditional UNet architecture [11], similar to the one used in [8], with modifications to make it suitable for EEG data. The DDPM's prediction is denoted as $\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)$.

### B. Conditional Autoencoder

The DDPM's forward pass leads to information loss, which the CAE attempts to make up for by recognizing and correcting these errors. This allows the CAE to generate more precise representations of the original EEG signals. To this end, we use the following objective function for the CAE:

$$\mathcal{L}_{\text{CAE}}(\psi, \phi) = ||\mathcal{L}_{\text{DDPM}}(\theta) - \mathcal{D}_\psi(\mathcal{E}_\phi(\mathbf{x}_0), \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t))||. \quad (5)$$

The CAE includes an encoder and decoder, denoted as $\mathcal{E}_\phi$ and $\mathcal{D}_\psi$, respectively. $\mathcal{D}_\psi$ is connected to the DDPM layers instead of the output of $\mathcal{E}_\phi$. This allows $\mathcal{D}_\psi$ to be implicitly conditioned on the corruption stage of the DDPM, as illustrated in Fig. 1 with *dashed arrows*. Additionally, to improve the reconstruction of $\mathcal{L}_{\text{DDPM}}$, the original signal, $\mathbf{x}_0$, and the output of the DDPM, $\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)$, are connected to the last layer of $\mathcal{D}_\psi$ as shown in Fig. 1 with *thin dashed arrows*.

### C. Classifier

After $\mathcal{E}_\phi$ has processed the data, the output is condensed into a single-dimensional representation, $\mathbf{z}$, using an adaptive average pooling layer. This creates a latent vector which is then fed into the linear classifier $\mathcal{C}_\rho$. The classifier is trained jointly with the CAE to differentiate the representations of each class and classify them. The dimension of $\mathbf{z}$ is fixed at 256 for the duration of the experiment. To include the classification loss in the CAE's objective function, we modified it to become the overall Diff-E objective.

$$\mathcal{L}_{\text{Diff-E}}(\psi, \phi, \rho) = ||\mathcal{L}_{\text{DDPM}}(\theta) - \mathcal{D}_\psi(\mathcal{E}_\phi(\mathbf{x}_0), \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t))|| \\ + \alpha||\hat{y} - y||_2. \quad (6)$$

The predicted label of the input signal is calculated as $\hat{y} = \mathcal{C}_\rho(\mathbf{z})$, where $\rho$ is an adjustable parameter for $\mathcal{C}_\rho$ and $y$ is the true label. The hyperparameter $\alpha$ is used to control the relative importance of the reconstruction loss and the classification loss, with a value of $0.1$ chosen for the experiment. During inference, only $\mathcal{E}_\phi$ and $\mathcal{C}_\rho$ are used to classify the signals, with the predicted label obtained as $\hat{y} = \mathcal{E}_\phi(\mathcal{C}_\rho(\mathbf{x}_0))$. To evaluate the effectiveness of Diff-E, it is compared to other methods that have been applied to decoding EEG signals in various paradigms, such as motor imagery and event-related potentials [12], [13]. This comparison is conducted to assess the performance of Diff-E and to determine its suitability for imagined speech EEG signal decoding applications.

### D. Model Implementation Details

In our study, the DDPM and CAE frameworks are constructed with layers that sequentially execute convolution, normalization, and activation functions. The encoder employs adaptive pooling to produce a compact feature vector, $\mathbf{z}$. The total number of trainable parameters for DDPM and CAE is roughly $3e+5$ and the classifier has $4e+5$ parameters. Optimization is conducted using RMSProp and a cyclic learning rate that starts at $9e-5$ and caps at $1.5e-3$. The training extended over 500 epochs, with L1 loss for DDPM and CAE, and mean squared error for the classifier's one-hot encoded classification tasks. For model evaluation, 20 % of the data was reserved for testing, with a consistent random seed ensuring reproducibility.

### E. Dataset

*1) Data Description:* This study used data from a previous study by Lee et al. [14], [15]. Participants were 22 healthy individuals, 15 of whom were male, with an average age of 24.68 ± 2.15. None of them had a history of neurological disease or language disorders and had no hearing or visual impairments. Furthermore, they did not take drugs for 12 hours prior to the session. All of them had received high-quality English education for more than 15 years. The overt speech task involved instructing the 22 subjects to imagine saying 12 different words or sentences, such as "ambulance," "clock," "hello," "help me," "light", "pain," "stop," "thank you," "toilet," "TV," "water" and "yes," as well as a resting state, resulting in a total of 13 classes. The researchers used a 64-channel EEG cap with active Ag/AgCl electrodes that followed the international 10-10 system to record EEG signals. The FCz and FPz channels were set as the reference and ground electrodes, respectively. Brain Vision/Recorder software (BrainProduct GmbH, Germany) was used to collect the EEG signals, which were then operated using the MATLAB 2018a software. The researchers made sure that the impedance of all electrodes was kept below 10 $k\Omega$. The researchers randomly presented 22 blocks of 12 words and a rest class. Each of the 22 participants contributed 1,300 samples, consisting of 100 samples per category. The study was approved by the Institutional Review Board of the Korea University [KUIRB-2019-0143-01] and was conducted in accordance with the Declaration of Helsinki.

*2) Preprocessing:* This research used a variety of preprocessing techniques to ensure the accuracy of the EEG data. Initially, a bandpass filter was used to filter signals between 0.5 and 125 Hz, with additional notch filtering at 60 and 120 Hz to eliminate power line interference. Subsequently, a common average reference method was used to reference the data and reduce any noise present. To remove ocular and muscular artifacts caused by movement or sounds, automatic electrooculography and electromyography removal methods were employed. After the artifacts were removed, the EEG signals were chosen in the high-gamma frequency band to train the model and analyze the data. The data set was then epoched into 2-second segments, with a baseline correction applied 500 ms before the task. All preprocessing steps were performed using MATLAB-based tools, such as the OpenBMI Toolbox [16], [17] or BBCI Toolbox [14].

### III. RESULTS AND DISCUSSION

In this study, we compared the performance of our proposed method with three established approaches: DeepConvNet [5], EEGNet [18], and the method introduced by Lee et al. [15], in the context of the decoding of the EEG signal from spoken speech. The results, presented in Table 1, show that our method outperformed the other three in terms of both accuracy and area under the curve (AUC). The average accuracy of our approach was 72.33 %, with a standard deviation of 7.51 %, while the average AUC was 93.22 %, with a standard deviation of 3.18 %. These figures are significantly higher than those of the compared methods. Specifically, the approach of DeepConvNet, EEGNet and Lee et al. yielded average average accuracies of 32.34 %, 42.73 %, and 57.06 %, and average AUCs of 73.00 %, 81.00 %, and 83.01 %, respectively.

TABLE I
ACCURACY AND AUC SCORES FOR IMAGINED SPEECH CLASSIFICATION

| Subject | Accuracy (%) | AUC (%) |
|---|---|---|
| DeepConvNet | 32.34 ± 5.10 | 73.00 ± 4.00 |
| EEGNet | 42.73 ± 3.80 | 81.00 ± 4.19 |
| Lee et al. | 57.06 ± 6.52 | 83.01 ± 5.10 |
| Diff-E | **72.33 ± 7.51** | **93.22 ± 3.18** |

TABLE II
EFFICACY OF EACH COMPONENT IN DIFF-E: AN ABLATION STUDY
ASSESSING THE INDIVIDUAL CONTRIBUTIONS OF DDPM AND CAE

| Components | Accuracy (%) | AUC (%) |
|---|---|---|
| Diff-E | **72.33 ± 7.51** | **93.22 ± 3.18** |
| w/o DDPM | 52.11 ± 8.98 | 90.19 ± 5.11 |
| w/o DDPM & $\mathcal{D}_\psi$ | 51.11 ± 8.80 | 66.53 ± 4.54 |

This indicates the superior ability of our proposed method in decoding EEG signals related to spoken speech.

Our research yielded unexpected results, especially since the more traditional approach of Lee et al. [15], which combines a common spatial pattern with the support vector machine, outperformed popular EEG decoding methods such as EEGNet [18] and DeepConvNet [5]. These methods have been extensively used for motor imagery and event-related potentials. Our findings emphasize the importance of selecting a suitable model architecture that is compatible with the task, the EEG paradigms used, and other relevant factors.

## IV. CONCLUSION

This research marks a significant advancement in the field of EEG signal decoding, particularly with regard to the challenge of interpreting spoken language. Our study introduces an innovative application of generative models that showcases improved performance over more conventional neural network approaches such as DeepConvNet and EEGNet. The findings suggest that generative models could be instrumental in improving the processing of EEG signals and could potentially be adapted for wider applications within this scientific area.

Furthermore, our research underscores the critical role of model architecture selection in EEG decoding tasks. The observed disparities in the performance of established methods like DeepConvNet and EEGNet, when compared to our generative model approach, underscore this point. The choice of model must be informed by a thorough understanding of the characteristics of the EEG data and the specific requirements of the decoding task.

In summary, our research offers a promising avenue for the accurate interpretation of EEG signals related to spoken language. This has far-reaching implications for the progression of brain-computer interfaces, potentially enhancing communication capabilities and assistive technologies. Additionally, it furthers the knowledge and application of deep learning in EEG analysis, potentially setting a precedent for future research in this vital area.

## REFERENCES

[1] K.-H. Thung *et al.*, "Conversion and time-to-conversion predictions of mild cognitive impairment using low-rank affinity pursuit denoising and matrix completion," *Med. Image Anal.*, vol. 45, pp. 68–82, 2018.

[2] J. Kim *et al.*, "Abstract representations of associated emotions in the human brain," *J. Neurosci.*, vol. 35, no. 14, pp. 5655–5663, 2015.

[3] Z. Tayeb *et al.*, "Validating deep neural networks for online decoding of motor imagery movements from EEG signals," *Sensors*, vol. 19, no. 1, p. 210, 2019.

[4] M. Lee, C.-B. Song, G.-H. Shin, and S.-W. Lee, "Possible effect of binaural beat combined with autonomous sensory meridian response for inducing sleep," *Front. Hum. Neurosci.*, vol. 13, pp. 425–440, 2019.

[5] R. T. Schirrmeister *et al.*, "Deep learning with convolutional neural networks for eeg decoding and visualization," *Hum. brain mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.

[6] K.-T. Kim, C. Guan, and S.-W. Lee, "A subject-transfer framework based on single-trial EMG analysis using convolutional neural networks," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 28, no. 1, pp. 94–103, 2019.

[7] R. Mane *et al.*, "FBCNet: A multi-view convolutional neural network for brain-computer interface," *arXiv preprint arXiv:2104.01233*, 2021.

[8] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *Adv. Neural Inf. Process. Syst. (NeurIPS)*, vol. 33, pp. 6840–6851.

[9] S.-H. Lee, M. Lee, J.-H. Jeong, and S.-W. Lee, "Towards an EEG-based intuitive BCI communication system using imagined speech and visual imagery," in *Conf. Proc. IEEE Int. Conf. Syst. Man Cybern. (SMC)*, pp. 4409–4414.

[10] H.-I. Suk, S. Fazli, J. Mehnert, K.-R. Müller, and S.-W. Lee, "Predicting BCI subject performance using probabilistic spatio-temporal filters," *PLoS One*, vol. 9, no. 2, p. e87056, 2014.

[11] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Med. Image Comput. Comput. Assist. Interv. (MICCAI)*, pp. 234–241.

[12] J.-S. Bang, M.-H. Lee, S. Fazil, C. Guan, and S.-W. Lee, "Spatio-spectral feature representation for motor imagery classification using convolutional neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 7, pp. 3038–3049, 2021.

[13] S.-B. Lee *et al.*, "Comparative analysis of features extracted from EEG spatial, spectral and temporal domains for binary and multiclass motor imagery classification," *Inf. Sci.*, vol. 502, pp. 190–200, 2019.

[14] R. Krepki, B. Blankertz, G. Curio, and K.-R. Müller, "The Berlin brain-computer interface (BBCI)–towards a new communication channel for online control in gaming applications," *Multimed. Tools Appl.*, vol. 33, no. 1, pp. 73–90, 2007.

[15] S.-H. Lee, M. Lee, and S.-W. Lee, "Neural decoding of imagined speech and visual imagery as intuitive paradigms for BCI communication," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 28, no. 12, pp. 2647–2659, 2020.

[16] M.-H. Lee *et al.*, "EEG dataset and OpenBMI toolbox for three BCI paradigms: An investigation into BCI illiteracy," *GigaScience*, vol. 8, no. 5, p. giz002, 2019.

[17] Y.-E. Lee and S.-W. Lee, "Decoding event-related potential from ear-EEG signals based on ensemble convolutional neural networks in ambulatory environment," in *Int. Winter Conf. Brain Comput. Interface (BCI)*, pp. 1–5.

[18] V. J. Lawhern *et al.*, "EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces," *J. neural engineering*, vol. 15, no. 5, p. 056013, 2018.