# Predicting EEG Responses to Attended Speech via Deep Neural Networks for Speech

Emina Alickovic, Tobias Dorszewski, Thomas U. Christiansen, Kasper Eskelund,
Leonardo Gizzi, Martin A. Skoglund, Dorothea Wendt

*Abstract*— Attending to the speech stream of interest in multi-talker environments can be a challenging task, particularly for listeners with hearing impairment. Research suggests that neural responses assessed with electroencephalography (EEG) are modulated by listener's auditory attention, revealing selective neural tracking (NT) of the attended speech. NT methods mostly rely on hand-engineered acoustic and linguistic speech features to predict the neural response. Only recently, deep neural network (DNN) models without specific linguistic information have been used to extract speech features for NT, demonstrating that speech features in hierarchical DNN layers can predict neural responses throughout the auditory pathway. In this study, we go one step further to investigate the suitability of similar DNN models for speech to predict neural responses to competing speech observed in EEG. We recorded EEG data using a 64-channel acquisition system from 17 listeners with normal hearing instructed to attend to one of two competing talkers. Our data revealed that EEG responses are significantly better predicted by DNN-extracted speech features than by hand-engineered acoustic features. Furthermore, analysis of hierarchical DNN layers showed that early layers yielded the highest predictions. Moreover, we found a significant increase in auditory attention classification accuracies with the use of DNN-extracted speech features over the use of hand-engineered acoustic features. These findings open a new avenue for development of new NT measures to evaluate and further advance hearing technology.

## I. INTRODUCTION

The challenge of attending the target speech signal while ignoring other sounds is an extensively-studied problem known as the "cocktail party" problem [2]. To date, research on "cocktail party" environments has progressed significantly to a point where it is now possible to decode (i.e., classify) the attended speech by quantifying "neural tracking" (NT) of speech [3], [4], i.e., by comparing the neural activity of the listener, recorded with electroencephalogram (EEG), to the activity of multiple candidate speech sources in a listening environment [5], [6]. NT methods have allowed to assess the real benefits of the hearing aids [6], [7].

NT methods involve encoding of the speech by estimating the temporal response function (TRF) that linearly maps time-lagged speech signals to EEG. NT methods proposed in the literature mostly rely on hand-engineered acoustic and linguistic speech feature that include for example speech envelope, spectrogram, pitch, phonetic and lexical features [8]–[12]. Despite the successful usage of acoustic-linguistic features in NT, it remains unclear to what extent language models [13], [14] may be used to extract relevant features.

Modern artificial intelligence models using deep neural networks (DNN) revolutionized the field of speech representation showing that DNN models, without specific language knowledge, can derive speech features that correlates well with neural responses recorded with functional magnetic resonance imaging (fMRI), magnetoencephalography (MEG) and electrocorticography (ECoG) [15]–[17]. Furthermore, recent work has suggested that DNN models can successfully predict speech comprehension from the neural responses recorded with fMRI [18]. Lastly, Li et. al. [17] demonstrated that speech features in hierarchical DNN layers can better predict neural responses than hand-engineered acoustic-phonetic features, as observed in high signal-to-noise (SNR) ECoG signals, throughout the auditory pathway [17]. However, whether DNN-derived speech features can better predict noninvasive, low SNR EEG responses and, classify attention in multi-talker listening environments remains unknown.

To address this issue, we compare a wide variety of hand-engineered acoustic and DNN-derived speech features in light of human neural responses to competing speech. Specifically, we analyze the EEG responses of 17 healthy younger adults. During 45 min-long sessions the listeners were instructed to attend to one of two competing talkers. We trained a variety of NT models and compare their ability to linearly map speech onto the EEG recordings. Furthermore, we investigate the hierarchy of layers in the DNN models. Finally, we trained a variety of NT models and compare their ability classify the attended speech.

## II. METHODS

The experimental protocol was reviewed and approved by the ethics committee for the capital region of Denmark (journal number H-21065001). The study was conducted according to the Declaration of Helsinki, and all the participants gave a written consent prior to the experiment.

### A. Study Design

*1) Participants:* Participants comprised 17 younger adults (9 males, mean age 29.0, SD 6.4). All participants were

native Danish speakers, had normal or corrected-to-normal vision, had no history of neurological disorders, dyslexia, or diabetes mellitus, and had clinically normal hearing.

*2) Stimuli and Recording:* During the experimental session presented in this study, participants were asked to attend to one of two different audiobooks narrated in Danish. The stimuli comprised 33 ∼1 min-long segments from audiobook recordings of *Himalaya i sigte* (a story of traveling in the Himalayas) read by a female talker and *Simon* (a biography on Simon Spies) read by a male talker, and sampled at 44.1 kHz. Prolonged silent periods in the speech stimuli longer than 500 ms were shortened to 500 ms. Stimuli were routed through a sound card (RME Hammerfall DSP multiface II, Audio AG, Germany) and were played via loudspeakers (Genelec 8040A; Genelec Oy, Finland) at an average intensity of 70 dB SPL each positioned ±30° to the left or right of the center. EEG data were acquired at a sampling rate of 8192 Hz with a BioSemi ActiveTwo 64-channel EEG recording system in 10-20 layout.

*3) Experimental Session Design:* A total of 33 trials were conducted, with the first trial used for training and 32 trials used for analysis. Each trial consisted of 5 s of silence, 1 min of speech mixture and two 2-choice questions on attended story to keep the participants alert. The 32 trials were divided into 8 blocks of 4 randomized consecutive trials, with 2 blocks for each of "male right", "male left", "female right" and "female left". Before each block, a visual cue on the screen and 5 second of the to-be-attended speech were presented indicating talker (male or female) and the side (left or right) to be attended.

*B. Neural Data Analysis*

*1) EEG Preprocessing:* The EEG data were re-referenced to the average of the mastoid electrodes, band-pass filtered between 0.1 and 10 Hz and re-sampled to 100 Hz. Subsequently, signal components of non-neural origin were removed using a procedure based on independent component analysis [19]. Next, the data were filtered between 1 and 10 Hz and normalized to zero mean and unit variance. In a final step, data were segmented into trials of 59 s duration from 0 to 59 s relative to the onset of the speech.

*2) Neural Tracking of Speech:*

*a) Quantifying Brain Prediction Scores:* The TRF framework allows to study how the brain processes competing speech. It includes two stages: a training stage to derive TRFs for the each talker and a testing stage to quantify how well EEG responses can be predicted. In the training stage, time-lagged speech features of each talker $S$ are linearly mapped to the EEG response(s) $R$ of the listener based on TRF $W$ derived via regularized linear regression (rLR) with a parameter $\lambda$ to control for overfitting [3]:

$$W = \left[\mathcal{H}(S)^T \mathcal{H}(S) + \lambda I\right]^{-1} \mathcal{H}(S)R \qquad (1)$$

where $\mathcal{H}(*)$ is a Hankel matrix (see [3] for more details). In the testing stage, EEG responses are predicted as:

$$\widehat{R} = \mathcal{H}(S)W \qquad (2)$$

The quality of the prediction is quantified in terms of a brain prediction score (BPS) measuring the correlation (a Pearson's $r$) between the true and reconstructed EEG responses.

*b) Classifying Auditory Attention:* The NT framework allows to classify auditory attention (i.e., identify the attended speech) in multi-talker environments. In order to classify which of the streams a listener attended to, two TRF models $W_{att}$ and $W_{ign}$ from (1) are assembled to become two competing prediction models for every single EEG channel. Next, the EEG signals ($\widehat{R}_{att}$ and $\widehat{R}_{ign}$) are independently predicted from the attended ($S_{att}$) and the ignored ($S_{ign}$) speech signals. To estimate which of the predicted EEG signals ($\widehat{R}_{att}$ versus $\widehat{R}_{ign}$) are most likely representing the attended speech, we compute channel-by-channel brain prediction scores $BPS_{att}$ and $BPS_{ign}$. Finally, we compare $BPS_{att}$ and $BPS_{ign}$ values averaged across all EEG channels and the signal with the highest BPS is classified as the attended speech.

*C. Speech Feature Extraction*

*1) Hand-Engineered Speech Features:* Hand-engineered acoustic features considered for this study included *speech envelope* (the root-mean-square of 10 ms windows and scaled by raising the value to the power of a compression parameter of 0.3 or 1 indicating no compression), the *envelope derivative*, *spectrogram* (100 linearly spaced components between 0 Hz and 8 kHz computed with a short-time Fourier transform at 100 Hz, and scaled by a compression parameter of 0.3), Mel Frequency Cepstral Coefficients (*MFCC*) (13 components representing different frequencies between 20 Hz and 8 kHz, as well as their derivatives and the second derivative for a total of 39 features) and *pitch* (absolute and relative pitch, and pitch change computed as in [17]). The MFCC feature set was also used as an initialization of the labels for DNN training. Similar to [17], we included a baseline model comprising full acoustic features.

*2) DNN-extracted Speech Features:* We employ the Hidden unit BERT (HuBERT) DNN model - a transformer-based self-supervised model for speech feature learning [14]. A major component of HuBERT model training is applying the predictive loss over the masked portions of speech driving the model to learn a fused acoustic and language feature set over the speech input. Using ECoG recordings, it has recently been shown that the HuBERT DNN model yielded the best BPS among the benchmarks DNN models [17].

Our goal is to extract relevant speech features in different DNN layers in order to use them as inputs to the NT models. With the speech material presented in our study being in Danish, different methods are employed to obtain a Danish version of HuBERT DNN model. For a complete description please refer to [1].

Since both the audio used in the experiments and the LibriSpeech corpus [20] of the English HuBERT model stem from audiobooks, we collected a similar Danish speech corpus for training purposes. Only continuous, clearly spoken, Danish speech without background noise was included. The main source was speech materials used in previous studies

| DNN Name | Data set | Weight source | Labels |
|----------|----------|---------------|--------|
| English | 960 h English | Random | 100, from MFCC |
| Danish | 65 h Danish | Random | 100, from MFCC |
| SSFT A1 | 65 h Danish | English | 250, from L9 of Danish |

conducted at Eriksholm Research Centre. Additionally, the publicly available Danish audiobooks found on LibriVox [21] were included. For all speech files, the starts of the files with an introduction (e.g., the name of the reader or the book title) were removed. Mono channel audio signal was created by averaging the stereo channels, resampled to $16\,kHz$ and divided into equally long segments. To avoid a strong effect of individual speakers, the amount of speech data was limited to 300 files per speaker. In total 3900 such audio files ($>$ $65\,h$) were prepared and divided into training and validation sets in a common $80/20$ split.

An important component of the HuBERT DNN models are the artificially created labels with which the DNN is trained. Based on MFCCs, all speech segments were clustered with k-means. These clusters were used as labels for the DNN training. We considered three major HuBERT DNN models, see Table I. First, the original 'base' HuBERT DNN model, referred to as the 'English DNN' in this study, was trained on $960\,h$ of continuous English speech from the LibriSpeech corpus [14]. Second, the Danish HuBERT DNN model, here referred to as the 'Danish DNN', used the $\sim 65\,h$ data set of Danish speech for training with randomly initialized weights which mimics the original training of the base English DNN [14]. Third, a self-supervised fine-tuning (SSFT) was added as an additional training option. The English DNN is used to initialize the network weights. Here, only the latest layers of the pre-trained DNN are replaced with new layers to allow prediction of different clusters. The training objective during SSFT is the HuBERT objective of classifying segments of the unseen speech data. We refer to this model as SSFT A1.

## III. RESULTS & DISCUSSION

### A. Predicting EEG from DNN-Extracted Speech Features

We first test whether DNN-extracted speech features linearly predict EEG responses. To this aim, we fit a rLR to predict the EEG activity elicited by the attended speech from the HuBERT model input with the same speech. We then compute a BPS, i.e. the correlation between the true EEG responses and the EEG responses predicted from the rLR. On average across EEG channels, the BPS for an English DNN HuBERT model are significantly distributed above zero with the mean BPS of 0.051 at layer 1 (L1), 0.054 at layer 5 (L5) and 0.05 at layer 12 (L12). Using speech features from the Danish DNN HuBERT model, a lower BPS of 0.035 at L1 to 0.015 at L12 were achieved.

Second, we evaluate whether DNN-extracted speech features can better predict EEG responses than hand-engineered acoustic features (see Fig. 1-2). Similar to the speech from HuBERT model, we first fit a rLR to predict the EEG

activity from the acoustic features. On average, a BPS of 0.039, 0.042, 0.029, 0.037, 0.032 and 0.038 for acoustic envelope, all envelope features, pitch, MFCC, spectrogram, and all features, respectively, were observed. Next, we compute a normalized BPS, i.e., the squared BPS from the DNN HuBERT model divided by the squared BPS with all envelope features as input to the rLR. Lastly, we observe that both the English DNN and the SSFT DNN A1 HuBERT models provide normalized BPS scores that are significantly higher than 1 at their best layer (p $<$ 0.05, 2-sided t-tests with BPS for 64 EEG channels). Our results are consistent with previous findings suggesting that DNN-derived speech features correlates well with the neural activity [15]–[17].



Fig. 1. The normalized mean brain prediction scores (BPS) for predicting the EEG activity elicited by the attended speech from the HuBERT model input with the same speech. A BPS higher than 1 signifies a higher quality of prediction from DNN-extracted speech features than from the hand-engineered acoustic features.

Third, analysis of hierarchical DNN layers shows that early layers (layer group 1: L1-L5) yielded higher BPS than later layers (layer group 2: L6-L12; $p < 0.05$ for one-way ANOVA factor 'layer group'). This is in line with recent studies providing evidence for the speech processing hierarchy within auditory cortex [22], [23].



Fig. 2. The mean BPS for each subject predicting the EEG responses from the HuBERT DNN models averaged across all trials and channels.

### B. Attention Classification with DNN-extracted Speech

Classifying auditory attention is notoriously challenging [4]. This issue poses strong limitations on the future application of NT methods to hearing devices. While hand-engineered acoustic features can be used in NT methods

to decode attention, we show that DNN-extracted speech features yield results that are consistently higher than those described in neuroscientific literature (see Fig. 3). We find a significant increase in classification accuracy with the use of DNN-extracted features over the use of acoustic features.

For the best acoustic feature set (envelope features), a mean attention (attended vs. ignored speech) classification accuracy of 75% (SD 43%) was achieved. With the best DNN feature set (from the L5 of the English DNN), the attention classification accuracy improved to 79% (SD 40%), yielding statistically significant differences (p = 0.0319, 2-sided t-test with 17 subjects).



Fig. 3. Auditory attention classification results per subject shown for the predictions based on the best acoustic feature set and the best DNN feature set (features from the layer 5 of an DNN trained only on English audio).

## IV. CONCLUSIONS

We propose a new framework to predict EEG responses to attended speech. Overall, the present study suggests DNN models for speech can retrieve information that correlate to speech processing hierarchy. Interestingly, our analyses highlights that EEG responses are significantly better predicted by DNN-extracted speech features than by hand-engineered acoustic features. Furthermore, analysis of hierarchical DNN layers shows that early layers yield the highest predictions. Finally, we find a significant increase in auditory attention classification accuracy with the use of DNN-extracted speech features over the use of hand-engineered acoustic features. In sum, NT methods could used to evaluate and further advance hearing technology and we propose a new approach to increase EEG prediction and attention classification accuracy.

REFERENCES

[1] T. Dorszewski, "Modeling brain tracking of speech with a deep neural network," University of Stuttgart, 2023.
[2] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears," *The Journal of the acoustical society of America*, vol. 25, no. 5, pp. 975–979, 1953.
[3] E. Alickovic, T. Lunner, F. Gustafsson, and L. Ljung, "A tutorial on auditory attention identification methods," *Frontiers in neuroscience*, p. 153, 2019.
[4] S. Geirnaert, S. Vandecappelle, E. Alickovic, A. de Cheveigne, E. Lalor, B. T. Meyer, S. Miran, T. Francart, and A. Bertrand, "Electroencephalography-based auditory attention decoding: Toward neurosteered hearing devices," *IEEE Signal Processing Magazine*, vol. 38, no. 4, pp. 89–102, 2021.
[5] J. A. O'Sullivan, A. J. Power, N. Mesgarani, S. Rajaram, J. J. Foxe, B. G. Shinn-Cunningham, M. Slaney, S. A. Shamma, and E. C. Lalor, "Attentional selection in a cocktail party environment can be decoded from single-trial EEG," *Cerebral cortex*, vol. 25, no. 7, pp. 1697–1706, 2015.
[6] E. Alickovic, T. Lunner, D. Wendt, L. Fiedler, R. Hietkamp, E. H. N. Ng, and C. Graversen, "Neural representation enhanced for speech and reduced for background noise with a hearing aid noise reduction scheme during a selective attention task," *Frontiers in neuroscience*, vol. 14, p. 846, 2020.
[7] E. Alickovic, E. H. N. Ng, L. Fiedler, S. Santurette, H. Innes-Brown, and C. Graversen, "Effects of hearing aid noise reduction on early and late cortical representations of competing talkers in noise," *Frontiers in Neuroscience*, vol. 15, p. 636060, 2021.
[8] W. Biesmans, N. Das, T. Francart, and A. Bertrand, "Auditory-inspired speech envelope extraction methods for improved EEG-based auditory attention detection in a cocktail party scenario," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 5, pp. 402–412, 2016.
[9] G. M. Di Liberto, J. A. O'sullivan, and E. C. Lalor, "Low-frequency cortical entrainment to speech reflects phoneme-level processing," *Current Biology*, vol. 25, no. 19, pp. 2457–2465, 2015.
[10] M. P. Broderick, A. J. Anderson, G. M. Di Liberto, M. J. Crosse, and E. C. Lalor, "Electrophysiological correlates of semantic dissimilarity reflect the comprehension of natural, narrative speech," *Current Biology*, vol. 28, no. 5, pp. 803–809, 2018.
[11] M. Gillis, J. Vanthornhout, J. Z. Simon, T. Francart, and C. Brodbeck, "Neural markers of speech comprehension: measuring EEG tracking of linguistic speech representations, controlling the speech acoustics," *Journal of Neuroscience*, vol. 41, no. 50, pp. 10 316–10 329, 2021.
[12] P. Patel, K. van der Heijden, S. Bickel, J. L. Herrero, A. D. Mehta, and N. Mesgarani, "Interaction of bottom-up and top-down neural mechanisms in spatial multi-talker speech perception," *Current Biology*, vol. 32, no. 18, pp. 3971–3986, 2022.
[13] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.
[14] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
[15] A. Goldstein, Z. Zada, E. Buchnik, M. Schain, A. Price, B. Aubrey, S. A. Nastase, A. Feder, D. Emanuel, A. Cohen *et al.*, "Shared computational principles for language processing in humans and deep language models," *Nature neuroscience*, vol. 25, pp. 369–380, 2022.
[16] C. Caucheteux and J.-R. King, "Brains and algorithms partially converge in natural language processing," *Communications biology*, vol. 5, no. 1, pp. 1–10, 2022.
[17] Y. Li, G. K. Anumanchipalli, A. Mohamed, J. Lu, J. Wu, and E. F. Chang, "Dissecting neural computations of the human auditory pathway using deep neural networks for speech," *bioRxiv*, 2022.
[18] C. Caucheteux, A. Gramfort, and J.-R. King, "Deep language algorithms predict semantic comprehension from brain activity," *Scientific reports*, vol. 12, no. 1, pp. 1–10, 2022.
[19] A. J. Bell and T. J. Sejnowski, "An information-maximization approach to blind separation and blind deconvolution," *Neural computation*, vol. 7, no. 6, pp. 1129–1159, 1995.
[20] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
[21] J. Kearns, "Librivox: Free public domain audiobooks," *Reference Reviews*, 2014.
[22] A. J. Kell, D. L. Yamins, E. N. Shook, S. V. Norman-Haignere, and J. H. McDermott, "A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy," *Neuron*, vol. 98, no. 3, pp. 630–644, 2018.
[23] J. Millet and J.-R. King, "Inductive biases, pretraining and fine-tuning jointly account for brain responses to speech," *arXiv preprint arXiv:2103.01032*, 2021.