



Extracting Different Levels of Speech Information from EEG Using an LSTM-Based Model

Mohammad Jalilpour Monesi^{1,2}, Bernd Accou^{1,2}, Tom Francart², Hugo Van hamme¹

¹KU Leuven, PSI, Dept. of Electrical engineering (ESAT), Leuven, Belgium

²KU Leuven, ExpORL, Dept. Neurosciences, Leuven, Belgium

mohammad.jalilpourmonesi@esat.kuleuven.be, hugo.vanhamme@esat.kuleuven.be

Abstract

Decoding the speech signal that a person is listening to from the human brain via electroencephalography (EEG) can help us understand how our auditory system works. Linear models have been used to reconstruct the EEG from speech or vice versa. Recently, Artificial Neural Networks (ANNs) such as Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) based architectures have outperformed linear models in modeling the relation between EEG and speech. Before attempting to use these models in real-world applications such as hearing tests or (second) language comprehension assessment we need to know what level of speech information is being utilized by these models. In this study, we aim to analyze the performance of an LSTM-based model using different levels of speech features. The task of the model is to determine which of two given speech segments is matched with the recorded EEG. We used low- and high-level speech features including: envelope, mel spectrogram, voice activity, phoneme identity, and word embedding. Our results suggest that the model exploits information about silences, intensity, and broad phonetic classes from the EEG. Furthermore, the mel spectrogram, which contains all this information, yields the highest accuracy (84%) among all the features.

Index Terms: LSTM, CNN, speech decoding, auditory system, EEG

1. Introduction

Understanding how natural running speech is processed in human auditory system has attracted a lot of attention in recent years. In the most common approach, natural running speech is presented to a listener and EEG signals are recorded simultaneously. A linear model is used to either reconstruct a speech representation from the recorded EEG (backward model) or to reconstruct the EEG from the speech (forward model). Then the correlation between the original and the reconstructed signal is used as a measure of neural speech tracking [e.g. 1, 2, 3, 4]. It has been shown that this approach can be used as an objective measure of how well speech is understood by a listener [3, 4, 5]. Real-time and accurate speech decoding from the brain has other potential applications such as Brain-Computer Interfaces (BCIs).

The speech envelope has been the most common choice as a representation of speech in this context [e.g. 6, 2, 3]. Nevertheless, low-level acoustic features such as the spectrogram [1, 7], or higher-level representations such as phonemes [1, 8], phonetic features [9, 7], phonotactics [10, 5], and semantics have also been used in linear backward/forward models. In [1, 11, 7], the authors have shown that including both low- and high-level speech features improves predicting brain responses

to speech. Lesenfants et al. [4] showed that combining phonetic features with the spectrogram improves speech reception threshold (SRT) prediction from EEG.

In the aforementioned studies linear models are used in a regression problem. Even though linear models are easy to interpret and implement, they are inherently not able to model complex and non-linear auditory processing in the brain. To address this, Artificial Neural Networks (ANNs) have been explored as an alternative approach. ANNs have been used in the context of Auditory Attention Decoding (AAD), in which the task is to determine which of two concurrent speakers a listener attends to using EEG signals. Simple feedforward neural networks [12] as well as Convolutional Neural Networks (CNNs) [13] have been used in AAD. In [14] we presented an LSTM-based model to relate EEG responses to speech stimuli through a match/mismatch classification problem. The match/mismatch classification task is defined here as which of two given speech stimuli caused a given EEG recording. We showed that the proposed model works better than the linear model and generalizes well to unseen subjects.

Our goal in this study is to understand what speech information we can infer from the recorded EEG using the LSTM-based model [14]. This will enable us to design speech decoding models for applications such as objective hearing tests and speech intelligibility measurement. Therefore, we analyze the performance of the LSTM-based model in the match/mismatch classification task using different levels of speech information. More specifically, we do a series of experiments where each uses a subset of speech information from the previous one and observe the difference in classification accuracy between these experiment. This enables us to see which speech information contributes to match/mismatch classification accuracy.

2. Methodology

In this section, we will first explain our data collection and preprocessing procedure. Then, we will introduce our classification task in detail. Finally, we present how we adjust the LSTM-based model [14] for new speech features such as the mel spectrogram.

2.1. Data collection and preprocessing

We recorded EEG signals from 86 normal-hearing native-Flemish-speaking subjects while they listened to natural running speech in the form of stories. Subjects were screened for normal hearing with pure tone audiometry and the Flemish matrix-test [15]. Throughout our recordings, stories were chosen from a set of 10 stories each lasting approximately 14 minutes and 30 seconds. 43 subjects listened to 8 stories, 38 listened to 7 stories, 4 listened to 6 stories, and 1 subject listened

to only 2 stories.

We used the APEX 4 software [16], developed at ExpORL, to present the stimuli (stories). The stimuli were presented binaurally at 62 dBA using Etymotic ER-3A insert phones. After each story, we asked subjects a question related to the story to make sure they paid attention. EEG signals were recorded using a 64 channel Active-two system from Biosemi at 8 kHz sampling rate. EEG recordings were measured inside an electromagnetically shielded and soundproofed cabin.

A multi-channel Wiener filter [17] was used to remove artifacts from the EEG recordings. Then, all EEG channels were re-referenced to common-average. Next, we bandpass filtered the EEG signals between 0.5-32 Hz using a Chebyshev2 filter with 80dB stopband attenuation and used Matlab's resample function to downsample them to 64 Hz. Matlab uses anti-aliasing low-pass filter before resampling. As a result the true bandwidth is cut to 0.5-30 Hz. Note that many studies [18, 3] have shown that Delta band (0.5-4 Hz) performs better than higher frequencies when relating speech to EEG. However, we included Theta, Alpha, and Beta frequency bands to provide more information to the model. As a last step, we applied mean-variance normalization for each single EEG recording, i.e., a recording of around 14 minutes and 30 seconds in response to one story.

Speech features: We used the following representations of speech:

1. Envelope: we used the 'powerlaw subbands' method [19] to extract the amplitude envelope, resulting in a one-dimensional feature.
2. Mel spectrogram: the mel scale was applied to Short Time Fourier Transform (STFT) in the range of 50 to 5000 Hz. We used 28 frequency bands.
3. Voice Activity Detection (VAD): This one-dimensional binary feature is one when the energy after pre-emphasis in a 15 ms frame is above the 75% percentile of values observed in the story, else it is zero.
4. Phoneme: Using forced alignment we segmented each word into a sequence of phonemes based on the International Phonetic Alphabet (IPA). Then for each of the 40 phonemes we converted the symbol to a one-hot vector of dimension 40.
5. Word embedding: We used pre-trained word embeddings of dimension 300 for Dutch which were trained using FastText [20].

The envelope and mel spectrogram were bandpass filtered between 0.5 and 32 Hz using a Chebyshev2 filter with 80dB stopband attenuation. All the speech features were downsampled to 64 Hz to be synchronous with the EEG.

2.2. Classification task

We split each recording into training, validation, and test sets. The validation and test sets contain 10% each, which were taken from the middle of the recording, and the remaining 80% were taken as the training set.

As the classification task is defined exactly the same for each of the speech features, we will only consider mel spectrogram (mel) here. We cut both EEG and mel into 5 seconds decision windows with 90% overlap. For each EEG segment of 5 seconds, we take the corresponding 5 seconds from the mel sequence as 'matched mel'. Next, we take a segment of mel that starts one second after the end of the matched mel as the 'mismatched mel' (see figure 1). Then, given (EEG, matched mel, mismatched mel) the task of the model is to determine which of the two mel spectrograms matches the EEG. Our outcome measure is the classification accuracy of this task.

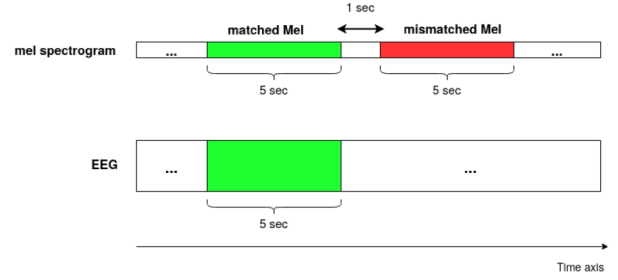


Figure 1: Extracting match and mismatch mel spectrograms with respect to an EEG segment of 5 seconds (decision window).

2.3. Model

In this study, we used the LSTM-based architecture described by [14], adjusted for different speech features. In figure 2 the LSTM-based model is shown with the mel spectrogram as a speech feature. The core ideas behind this architecture were (1) to use an LSTM layer in the speech path to model the brain response to the stimulus, including its delay, and (2) to map short segments of EEG and speech features to a common embedding space where we expect that matched speech (mel spectrogram) and EEG have similar representations while mismatched speech and EEG have dissimilar representations. For more details about the model and hyper-parameter tuning please see [14].

Here, compared to the original LSTM-based model, we changed the first layer in the EEG path (first 4 layers in figure 2 that apply to EEG) from a 2D convolution to a 1D convolution, for interpretation reasons. For some speech features we slightly adjusted the network: the EEG path remained the same but we added or removed some layers depending on the speech feature. In case of the mel spectrogram and phoneme-based features with feature dimension of at least 2, the architecture is shown in figure 2. For one-dimensional features such as Envelope and VAD, we do not have the first 1D convolutional layer in the speech path. For the word embedding feature, we replaced the 1D and 2D convolutional layers in the speech path with a max pooling layer with a stride of 3. For all the speech features, layers in the speech path are shared between the two speech inputs. We have provided the code for our models in <https://github.com/jalilpour-m/match-mismatch-speech-features>.

3. Results

For each speech feature, we trained the model in a subject independent fashion, i.e., we only trained one general model using all the training data of all the subjects. Also, when we concatenated two (or more) features, they are treated separately by the model up to the input of the LSTM layer, where for each time step we concatenate them in the feature space. We assume the LSTM is capable of integrating both information sources to the benefit of the match/mismatch classification accuracy.

3.1. VAD, envelope, and mel spectrogram

We wanted to know what level of speech information - reflected in EEG and speech features - the model exploits for the task. Therefore, we investigated how different speech features performed in our match/mismatch classification task. Previously

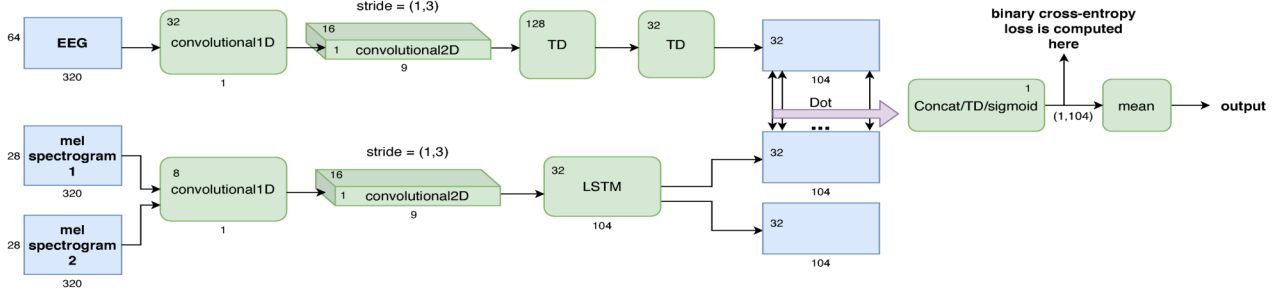


Figure 2: The LSTM-based model [14] for match/mismatch classification. TD refers to time distributed which applies a dense layer to every temporal slice of the input. Dot is a layer that applies dot product (cosine similarity) between EEG representation and speech representation for each time step.

we had evaluated the LSTM-based model working with the envelope as a speech feature. Here, we use the mel spectrogram which is a richer feature than the envelope (which can be computed from mel spectrogram), potentially leading to higher performance. In turn, from the envelope, we can derive the VAD feature. Figure 3 shows that the model’s median accuracy is 75% and 82% for VAD and envelope respectively. These results suggest that a substantial share of the accuracy comes from extracting silences from the EEG. Furthermore, when we concatenate VAD with envelope we do not gain any accuracy since VAD could be extracted from envelope. Moreover, concatenating mel spectrogram with envelope performs the same as mel spectrogram, as the mel spectrogram includes all the information available in the envelope. More specifically, there is a slight increase in performance when we use mel spectrogram instead of envelope, which is statistically significant ($z = -6.48$, $p < 0.001$). We used Wilcoxon signed-rank test with normal approximation for all our statistical tests.

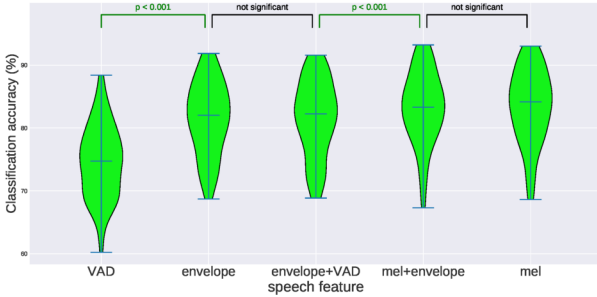


Figure 3: Classification accuracy of speech features voice activity detection (VAD), envelope, and mel spectrogram (referred to as mel). Violin plots are shown over 86 subjects. Decision window was set to 5 seconds.

3.2. Phoneme and word embedding

We also used a 40-dimensional phoneme feature as this has been shown to improve performance in linear models [7, 4]. We also created three simplified features from the phoneme feature where each feature contains less information compared to the previous one. First, we grouped phonemes into 5 categories: short vowels, long vowels, plosives, fricatives, nasals and approximants. We concatenated these 5 categories with silence and created a six-dimensional feature called broad phonetic classes (BPC) in this study. Second, we grouped vowel

els and consonants and together with silences creating a three-dimensional feature called vowel/consonant. Third, we combined vowels and consonants in the vowel/consonant feature to one category and with silences created a two-dimensional feature called anyPhoneme. AnyPhoneme is a one-hot vector which determines at each time step whether it is silence or any phoneme. The idea behind creating these simpler phoneme representations is to investigate what phonetic detail is used by the network. However, when we combine different phonemes into one category we lose phone onsets. This is especially the case for anyPhoneme feature which only retains onsets after pauses, which are only allowed between words by the forced alignment model. In the case of vowel/consonant, whenever there are two or more consecutive consonants we will not have onset information of the consecutive consonants. We will only have onsets when altering from vowel to consonant and vice versa. Similarly, we might lose some onsets in the BPC feature which will happen rarely. In order to ensure that performance of the model is not affected by losing onset information we also created all three features from phoneme onsets. The versions that we kept the phone onsets are named anyPhoneme_onset, vowel/consonant_onset, and BPC_onset in this study.

Performance of the model using phoneme related features is shown in figure 4. Note that we do not report accuracy for vowel/consonant_onset and BPC_onset since they had the same performance as vowel/consonant and BPC respectively. However, AnyPhoneme_onset outperforms anyPhoneme ($z = -8.01$, $p < 0.001$) suggesting that the model uses phone onset information. Next, we observe that vowel/consonant (or vowel/consonant_onset) has the same performance as anyPhoneme_onset. This suggests that the distinction between vowels and consonants is not informative enough. Nevertheless, it can use the phone onsets that occur in vowel/consonant to increase the accuracy compared to anyPhoneme feature. BPC outperforms vowel/consonant ($z = -5.37$, $p < 0.001$). This suggests that the model is capable of using some broad phonetic class information from the EEG. Finally, we observed no difference in performance when we used phoneme feature compared to the BPC. We argue that the model is unable to use more detailed phoneme information beyond BPC. Furthermore, concatenating mel with phoneme yields the same performance as that of the mel. This suggests that the model can extract phoneme (or BPC) level information from the mel spectrogram.

To test whether the model can extract any contextual information from speech, we used a word embedding feature. Figure 5 shows that the performance of the word embedding, referred to as wordEmb, is around only 70% median accuracy. Also,

when we concatenate this feature with envelope, we observe the same performance as for that of the envelope. In addition, the anyPhoneme feature which is basically a simplified binary version of wordEmb performs the same as wordEmb. These two observations suggest that the model does not use any word level information from the wordEmb feature.

Furthermore, we concatenated the envelope and the BPC (or phoneme) feature expecting that the performance will be close to that of the mel spectrogram. Figure 5 shows that concatenating the envelope with the BPC (env+BPC) or phoneme (env+phoneme) indeed increases the classification accuracy to almost that of the mel spectrogram. In addition, classification accuracy of combined envelope, phoneme, and mel spectrogram features is the same as for the mel spectrogram. These results are plausible since one could argue that the mel spectrogram contains intensity as well as BPC level information. If we average all the frequency bands of the mel spectrogram we obtain envelope which contains intensity (energy) of the speech. It has been shown a long time ago [21] that a simple multilayer perceptron with one layer and a few nodes can do BPC classification using mel spectrogram input. Note that the LSTM layer in our model is capable of behaving as a simple multilayer perceptron by always forgetting the past. In short, these observations suggest that within the significance of this experiment, the mel-model is exploiting intensity information as well as some information related to broad phonetic classes.

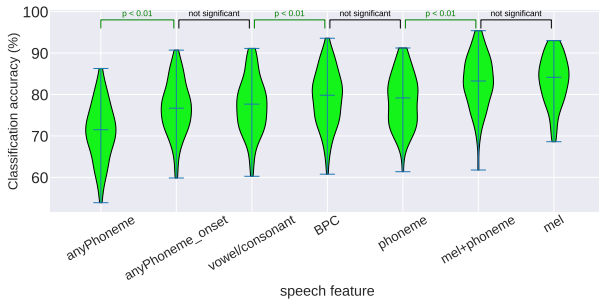


Figure 4: classification accuracy of the model for 5 phoneme related features and mel spectrogram. BPC stands for broad phonetic classes. Violin plots are shown over 86 subjects and the decision window is 5 seconds.

4. Conclusion

Our objective was to investigate what speech information we can decode from the recorded EEG. As a result, we trained the LSTM-based model [14] in the match/mismatch classification task using different levels of speech features. We used a broad range of speech features from a very low-level acoustic features such as VAD to higher-level features such as word embedding. In one set of experiments, we started with VAD because we wanted to test how much of the classification accuracy comes from just tracking silences in the speech. Our results show that indeed using only VAD the model can reach a classification accuracy of 75%. Then, we used envelope which contains intensity information over time. We observed that the classification accuracy increased to 82%. This suggest that the LSTM-based model is capable of extracting information about intensity of the speech. Finally, using mel spectrogram, which is a richer signal than the envelope, the model reached classification accuracy of 84%. The observed results are expected because envelope con-

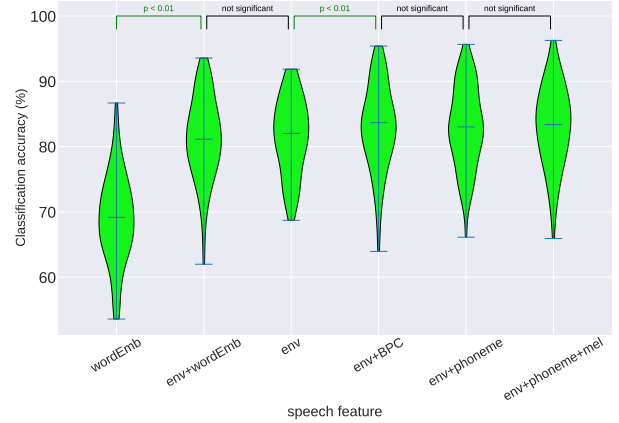


Figure 5: Classification accuracy of the model for word embedding (wordEmb), envelope, envelope+BPC, envelope+phoneme, and envelope+phoneme+mel spectrogram. In the figure, env, BPC, and mel refer to envelope, broad phonetic classes, and mel spectrogram respectively.

tains all the information present in VAD and in turn mel contains all the information present in the envelope.

In another set of experiments, we investigated which phoneme level information is being used by the model. We created three simpler phoneme-based features from the 40 dimensional phoneme feature: broad phonetic classes (BPC), vowel/consonant, and anyPhoneme. Our results show that phone onsets are important information for the performance of the model. More specifically, using the anyPhoneme feature we lose 5% accuracy compared to the anyPhoneme_onset feature. Considering phonetic detail, our results suggest that the model is able to use information up to the BPC feature. In other words, using the BPC feature or more detailed phoneme feature results in the same performance. We also showed that if we train the model with envelope and BPC features combined we reach the same performance as with that of the mel spectrogram. We argue that at least traces of BPC can be found in the EEG since if we remove this information from the speech path the classification accuracy decreases.

We also used a high level word embedding feature hoping that the model can extract some word level information from the EEG. We used this feature both alone and also in combination with the mel spectrogram. Unfortunately, our results suggest that the model is unable to extract word level information from EEG. In summary, it seems that the model is able to extract silences, energy intensity, and some broad phonetic class level information from EEG. Extracting word level information from EEG might be possible using more complex ANN models and more training data.

5. Acknowledgements

The work is funded by KU Leuven Special Research Fund C24/18/099 (C2 project to Tom Francart and Hugo Van hamme) and the Flemish Government under "Onderzoeksprogramma AI Vlaanderen". This project has also received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant agreement No 637424, ERC starting Grant to Tom Francart).

6. References

- [1] G. DiLiberto, J. OSullivan, and E. Lalor, "Low-frequency cortical entrainment to speech reflects phoneme-level processing," *Current Biology*, vol. 25, no. 19, pp. 2457–2465, 2015.
- [2] M. J. Crosse, G. M. Di Liberto, A. Bednar, and E. C. Lalor, "The Multivariate Temporal Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to Continuous Stimuli," *Frontiers in Human Neuroscience*, vol. 10, 2016.
- [3] J. Vanthornhout, L. Decruy, J. Wouters, J. Z. Simon, and T. Francart, "Speech Intelligibility Predicted from Neural Entrainment of the Speech Envelope," *Journal of the Association for Research in Otolaryngology*, vol. 19, no. 2, pp. 181–191, 2018.
- [4] D. Lesenfants, J. Vanthornhout, E. Verschueren, L. Decruy, and T. Francart, "Predicting individual speech intelligibility from the cortical tracking of acoustic- and phonetic-level speech representations," *Hearing Research*, vol. 380, pp. 1–9, 2019.
- [5] G. M. Di Liberto, J. Nie, J. Yeaton, B. Khalighinejad, S. A. Shamma, and N. Mesgarani, "Neural representation of linguistic feature hierarchy reflects second-language proficiency," *NeuroImage*, vol. 227, p. 117586, Feb. 2021.
- [6] N. Ding and J. Z. Simon, "Emergence of neural encoding of auditory objects while listening to competing speakers," *Proceedings of the National Academy of Sciences*, vol. 109, no. 29, pp. 11 854–11 859, Jul. 2012.
- [7] D. Lesenfants, J. Vanthornhout, E. Verschueren, and T. Francart, "Data-driven spatial filtering for improved measurement of cortical tracking of multiple representations of speech," *Journal Of Neural Engineering*, vol. 16, no. 6, Dec. 2019.
- [8] B. Khalighinejad, G. Cruzatto da Silva, and N. Mesgarani, "Dynamic Encoding of Acoustic Features in Neural Responses to Continuous Speech," *The Journal of Neuroscience*, vol. 37, no. 8, pp. 2176–2185, Feb. 2017.
- [9] E. S. Teoh and E. C. Lalor, "Attention differentially affects acoustic and phonetic feature encoding in a multispeaker environment," *bioRxiv*, p. 2020.06.08.141234, Jun. 2020, publisher: Cold Spring Harbor Laboratory Section: New Results. [Online]. Available: <https://www.biorxiv.org/content/10.1101/2020.06.08.141234v1>
- [10] G. M. Di Liberto, D. Wong, G. A. Melnik, and A. de Cheveign, "Low-frequency cortical responses to natural speech reflect probabilistic phonotactics," *NeuroImage*, vol. 196, pp. 237–247, Aug. 2019.
- [11] G. M. Di Liberto and E. C. Lalor, "Indexing cortical entrainment to natural speech at the phonemic level: Methodological considerations for applied research," *Hearing Research*, vol. 348, pp. 70–77, May 2017.
- [12] T. de Taillez, B. Kollmeier, and B. T. Meyer, "Machine learning for decoding listeners attention from electroencephalography evoked by continuous speech," *European Journal of Neuroscience*, vol. 0, no. 0, Dec. 2017.
- [13] G. Ciccarelli, M. Nolan, J. Perricone, P. T. Calamia, S. Haro, J. OSullivan, N. Mesgarani, T. F. Quatieri, and C. J. Smalt, "Comparison of Two-Talker Attention Decoding from EEG with Nonlinear Neural Networks and Linear Methods," *Scientific Reports*, vol. 9, no. 1, p. 11538, Dec. 2019.
- [14] M. J. Monesi, B. Accou, J. Montoya-Martinez, T. Francart, and H. V. Hamme, "An LSTM Based Architecture to Relate Speech Stimulus to Eeg," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2020, pp. 941–945.
- [15] H. Luts, S. Jansen, W. Dreschler, and J. Wouters, "Development and normative data for the Flemish/Dutch Matrix test," 2014. [Online]. Available: <https://lirias.kuleuven.be/retrieve/293640>
- [16] T. Francart, A. van Wieringen, and J. Wouters, "APEX 3: a multi-purpose test platform for auditory psychophysical experiments," *Journal of Neuroscience Methods*, vol. 172, no. 2, pp. 283–293, Jul. 2008.
- [17] B. Somers, T. Francart, and A. Bertrand, "A generic EEG artifact removal algorithm based on the multi-channel Wiener filter," *Journal of Neural Engineering*, vol. 15, no. 3, p. 036007, Feb. 2018.
- [18] N. Ding and J. Z. Simon, "Adaptive Temporal Encoding Leads to a Background-Insensitive Cortical Representation of Speech," *Journal of Neuroscience*, vol. 33, no. 13, pp. 5728–5735, Mar. 2013. [Online]. Available: <https://www.jneurosci.org/content/33/13/5728>
- [19] 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, May 2017.
- [20] E. Grave, P. Bojanowski, P. Gupta, A. Joulin, and T. Mikolov, "Learning word vectors for 157 languages," in *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [21] J. P. Martens and L. Depuydt, "Broad phonetic classification and segmentation of continuous speech by means of neural networks and dynamic programming," *Speech Communication*, vol. 10, no. 1, pp. 81–90, Feb. 1991.