

Speech Understanding Oppositely Affects Acoustic and Linguistic Neural Tracking in a Speech Rate Manipulation Paradigm

Eline Verschueren,^{1*} Marlies Gillis,^{1*} Lien Decruy,² Jonas Vanthornhout,¹ and Tom Francart¹

¹Research Group Experimental Oto-rhino-laryngology, Department of Neurosciences, KU Leuven-University of Leuven, Leuven, 3000, Belgium, and

²Institute for Systems Research, University of Maryland, College Park, Maryland 20742

When listening to continuous speech, the human brain can track features of the presented speech signal. It has been shown that neural tracking of acoustic features is a prerequisite for speech understanding and can predict speech understanding in controlled circumstances. However, the brain also tracks linguistic features of speech, which may be more directly related to speech understanding. We investigated acoustic and linguistic speech processing as a function of varying speech understanding by manipulating the speech rate. In this paradigm, acoustic and linguistic speech processing is affected simultaneously but in opposite directions: When the speech rate increases, more acoustic information per second is present. In contrast, the tracking of linguistic information becomes more challenging when speech is less intelligible at higher speech rates. We measured the EEG of 18 participants (4 male) who listened to speech at various speech rates. As expected and confirmed by the behavioral results, speech understanding decreased with increasing speech rate. Accordingly, linguistic neural tracking decreased with increasing speech rate, but acoustic neural tracking increased. This indicates that neural tracking of linguistic representations can capture the gradual effect of decreasing speech understanding. In addition, increased acoustic neural tracking does not necessarily imply better speech understanding. This suggests that, although more challenging to measure because of the low signal-to-noise ratio, linguistic neural tracking may be a more direct predictor of speech understanding.

Key words: acoustic hearing; linguistic representations; natural speech; neural coding

Significance Statement

An increasingly popular method to investigate neural speech processing is to measure neural tracking. Although much research has been done on how the brain tracks acoustic speech features, linguistic speech features have received less attention. In this study, we disentangled acoustic and linguistic characteristics of neural speech tracking via manipulating the speech rate. A proper way of objectively measuring auditory and language processing paves the way toward clinical applications: An objective measure of speech understanding would allow for behavioral-free evaluation of speech understanding, which allows to evaluate hearing loss and adjust hearing aids based on brain responses. This objective measure would benefit populations from whom obtaining behavioral measures may be complex, such as young children or people with cognitive impairments.

Received Feb. 4, 2022; revised June 29, 2022; accepted July 17, 2022.

Author contributions: E.V., M.G., L.D., J.V., and T.F. designed research; E.V. and M.G. performed research; E.V., M.G., and L.D. analyzed data; E.V. wrote the first draft of the paper; E.V., M.G., L.D., J.V., and T.F. edited the paper; E.V. and M.G. wrote the paper.

This work was supported by European Research Council, European Union's Horizon 2020 research and innovation program (Grant Agreement 637424) to T.F.; and Research Foundation Flanders PhD Grant SB 1586118N to E.V., PhD Grant SB 15A0620N to M.G., and Postdoctoral Grant 1290821N to J.V. We thank Sofie Keunen and Elise Verwaerde for help in data acquisition.

*E.V. and M.G. contributed equally to this work as co-first authors.

The authors declare no competing financial interests.

Correspondence should be addressed to Eline Verschueren at eline.verschueren@kuleuven.be.

<https://doi.org/10.1523/JNEUROSCI.0259-22.2022>

Copyright © 2022 the authors

Introduction

Understanding speech relies on the integration of different acoustic and linguistic properties of the speech signal. The acoustic properties are mainly related to sound perception, while the linguistic properties are linked to the content and understanding of speech. When listening to continuous speech, our brain can track both the acoustic and linguistic properties of the presented speech signal.

Neural tracking of acoustic properties of natural speech has been the subject of many studies. Particular emphasis has been placed on recovering the temporal envelope, that is, the slow modulations of the speech signal, from the brain responses, so-called neural envelope tracking. The temporal envelope is essential for

speech understanding (Shannon et al., 1995), and neural envelope tracking can be linked to speech intelligibility (e.g., Ding and Simon, 2013; Vanthornhout et al., 2018; Lesenfants et al., 2019; Iotzov and Parra, 2019; Verschuereen et al., 2020). However, only taking acoustic speech properties into account to investigate neural speech tracking would underestimate the complexity of the human brain, where linguistic properties also play an essential part, as reviewed in detail by Brodbeck and Simon (2020).

In addition to acoustic properties, there is growing interest in retrieving linguistic properties from brain responses to speech. Broderick et al. (2018) used semantic dissimilarity to quantify the meaning carried by words based on their preceding context. They report that the brain responds in a time-locked way to the semantic context of each content word. Additionally, neural tracking is also observed to linguistic properties derived from the probability of a given word or phoneme (i.e., word or phoneme surprisal) (Brodbeck et al., 2018; Koskinen et al., 2020; Weissbart et al., 2020). Recently, Gillis et al. (2021) combined several linguistic neural tracking measures and evaluated the potential of each measure as a neural marker of speech intelligibility. After controlling for acoustic properties, phoneme surprisal, cohort entropy, word surprisal, and word frequency were significantly tracked. These results show the potential of linguistic representations as a neural marker of speech intelligibility. In addition, this underlines the importance of controlling for acoustic features when investigating linguistic neural processing, as acoustic and linguistic features are often correlated (Brodbeck and Simon, 2020).

We investigated whether neural speech processing can capture the effect of gradually decreasing speech understanding by manipulating the speech rate. In this study, we focused on acoustic and linguistic speech processing. By changing the speech rate, we manipulate acoustic and linguistic speech processing simultaneously but in opposite directions: When increasing the speech rate, more phonemes, words, and sentences, and thus more acoustic information per second are present. In contrast, linguistic tracking decreases because it becomes more challenging to identify the individual phonemes or words at high speech rates, causing decreased speech understanding. We hypothesize that neural tracking of acoustic features will increase with increasing speech rate because more acoustic information will be present. However, linguistic speech tracking will decrease with increasing speech rate because of decreasing speech understanding. The effect of speech rate on neural responses to speech has already been investigated. However, all these studies only investigated brain responses to the acoustic properties of the speech signal (Ahissar et al., 2001; Nourski et al., 2009; Hertrich et al., 2012; Müller et al., 2019; Casas et al., 2021). It is challenging to extrapolate these findings to what would happen at a linguistic level as acoustic aspects of the speech dominate the neural responses to speech. No study, to our knowledge, reported on how speech rate affects linguistic speech processing and the potential interaction between both. In addition, no consensus has been reached on the effect of speech rate on acoustic neural tracking. For example, Nourski et al. (2009) reported that phase-locked responses decrease with increasing speech rate, similar to Ahissar et al. (2001) and Hertrich et al. (2012). However, in the same data, Nourski et al. (2009) also reported that time-locked responses to the envelope (70–250 Hz) could still be found at very high speech rates where speech is no longer understood.

We investigated how linguistic and acoustic speech tracking is affected when speech understanding gradually decreases. Analyzing neural speech tracking to different characteristics of

the presented speech allows us to identify neural patterns associated with speech understanding.

Materials and Methods

Participants

Eighteen participants between 19 and 24 years of age (4 men and 14 women) took part in the experiment after having provided informed consent. Participants had Dutch as their mother tongue and were all normal-hearing, confirmed with pure tone audiometry (thresholds ≤ 25 dB HL at all octave frequencies from 125 Hz to 8 kHz). The study was approved by the Medical Ethics Committee UZ Leuven/Research (kU Leuven) with reference S57102.

Speech material. The story presented during the EEG measurement was “A casual vacancy” by J.K. Rowling, narrated in Dutch by Nelleke Noordervliet. The story was manually cut into 12 blocks of varying length selected from the following list: 4, 5, 8.5, 12.5, 18, and 23 min. After cutting the story, the story was time-compressed with the Pitch Synchronous Overlap and Add algorithm from PRAAT (Boursma and Weenink, 2018) to manipulate the speech rate. This algorithm compresses the speech by dividing the speech into small overlapping segments and then repeating or deleting segments to decrease or increase the speech rate, thereby preserving the pitch. However, this manipulation is not exactly comparable with fast natural speech where vowels and silences become shorter, while some consonants (like, e.g., ‘t’) cannot be shortened. Using this algorithm, all phonemes are merging similarly and the formant transitions are probably not completely realistic anymore. Six different compression ratios (CRs) were used: 1.4, 1.0, 0.6, 0.4, 0.28, 0.22 with corresponding speech rates varying from ≈ 2.6 syllables/s (CR = 1.4) to ≈ 16.2 syllables/s (CR = 0.22). An example audio fragment for every speech rate is added as Extended Data Figures 1-1, 1-2, 1-3, 1-4, 1-5, and 1-6. The fastest CR (CR = 0.22) was applied to the longest part (23 min), the one but fastest CR (CR = 0.28) to the one but longest part (18 min), and so on. This way, all story parts were compressed or expanded to ≈ 5 min. These blocks had slightly different lengths because word and sentence boundaries were taken into account while cutting the story, which is important for the linguistic analysis. Every speech rate was presented twice to obtain 10 min of speech at the same rate. The story was presented in chronological order, and no part of the story was presented twice to the same participant. Although the linguistic measures only go back up to four words, and higher top-down processing has limited impact on the linguistic measures we used, we presented all speech rates in random order per participant. This way, the effect of the level of understanding of previous parts, which could influence linguistic processing, was minimized. For each stimulus block, we determined the number of syllables using the forced aligner of the speech alignment component of the reading tutor (Duchateau et al., 2009) and CELEX database (Baayen et al., 1996). The number of syllables uttered for each speech block was then divided by the duration of the speech block in seconds to obtain the speech rate.

After each part of the story, content questions were asked to maximize the participants’ attention and motivation. In addition, speech intelligibility was measured after each block by asking the participants to rate their speech understanding on a scale from 0% to 100% following the question: “Which percentage of the story did you understand?” A short summary of the story was shown in the beginning of the experiment to enhance intelligibility as some participants started with more difficult speech rates.

Experimental setup

EEG recording. EEG was recorded with a 64-channel BioSemi ActiveTwo EEG recording system at a sample rate of 8192 Hz. Participants sat in a comfortable chair and were asked to move as little as possible during the EEG recordings. All stimuli were presented bilaterally using APEX 4 (Francart et al., 2008), an RME Multiface II sound card (Haimhausen), and Etymotic ER-3A insert phones. The setup was calibrated using a 2 cm³ coupler of the artificial ear (Brüel & Kjær 4152). Recordings were made in a sound-proof and electromagnetically shielded room.

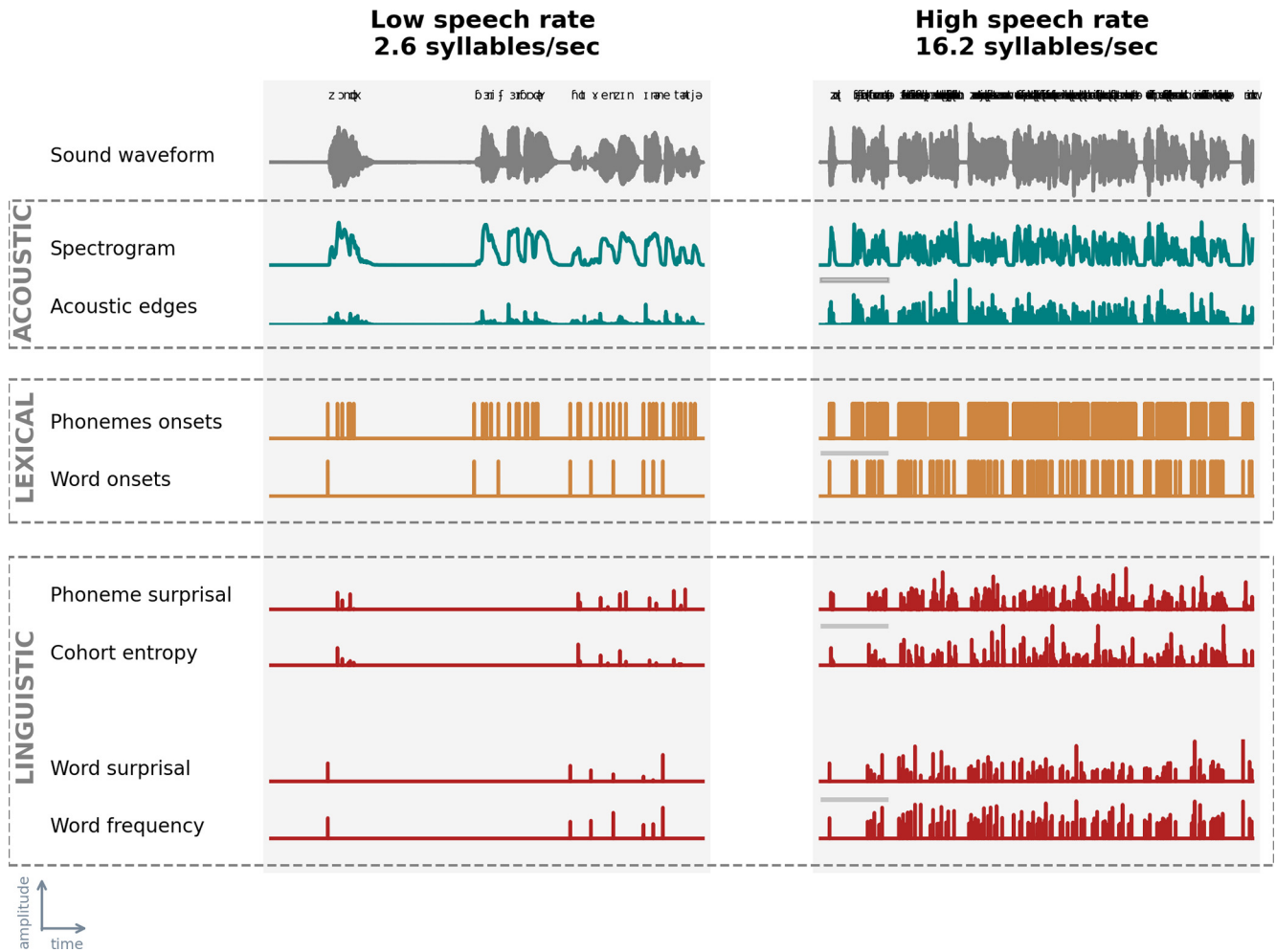


Figure 1. Speech representations at acoustic, lexical, and linguistic level. We visualized the speech representations used in this study for all three levels: acoustic (averaged across frequency bands; top box; blue), lexical (middle box; orange), and linguistic (bottom box; red) for the lowest speech rate (SR = 2.6 syllables/s, left) and highest speech rate (SR = 16.2 syllables/s; right) for the first 10 s of the speech material. This is the first sentence of the audio book “Barry Fairbrother” (second and third word) is a proper name and therefore not included in the corpus used for the linguistic features.

Signal processing

EEG processing. We processed the EEG in 5 consecutive steps. First, we drift-corrected the EEG signals by applying a first-order highpass Butterworth filter with a cutoff frequency of 0.5 Hz in the forward and backward direction. Then, we reduced the sampling frequency of the EEG from 8192 to 256 Hz to reduce computation time. Artifacts related to eyeblinks were removed with a multichannel Wiener filter (Somers et al., 2018). Subsequently, we referenced the EEG signals to the common average signal. The major advantage of common average referencing is that topography patterns are not biased relative to the chosen channel. No particular accentuation of certain channels allows us to investigate the contribution of neural processes related to language processing evoked by spatially different neural sources. Last, we removed the power line frequency of 50 Hz by using a second-order IRR notch filter at this frequency with a quality factor of 35 to determine the filter’s bandwidth at -3 dB.

Stimuli representations. This study aims to investigate acoustic and linguistic neural tracking at different speech rates. To examine acoustic tracking, we estimated neural tracking based solely on acoustic representations of the stimulus, namely, the spectrogram and acoustic edges. To investigate linguistic neural tracking, we created two models: a model to control speech acoustics, which consisted of acoustic and lexical segmentation representations, and a model that included linguistic representations on top of these acoustic and lexical segmentation representations. All speech representations used in the analysis are visualized in Figure 1.

The spectrogram representations were calculated based on the low-pass filtered speech stimulus (zero-phase low-pass FIR filter with a hamming window of 159 samples). We low-pass filtered the stimulus at a cutoff frequency of 4 kHz as the insert earphones also low-pass filter at this frequency. Subsequently, we calculated the spectrogram representation from this filtered stimulus using the Gammatone Filterbank Toolkit (Heeris, 2014) (center frequencies between 70 and 4000 Hz with 256 filter channels and an integration window of 0.01 s). By using a Gammatone Filterbank, the estimated filter outputs are closer to the human auditory response: the higher the frequency, the wider the filter (Slaney, 1998). We combined the filter outputs by averaging them into eight frequency bands with center frequencies of 124, 262, 455, 723, 1098, 1618, 2343, and 3352 Hz. To calculate the acoustic edges representations, we took the derivative of the spectrogram’s response in each frequency band and mapped all its negative values to 0. Last, we reduced the sampling frequency of these representations to the same sampling frequency as the EEG (i.e., 256 Hz).

To determine linguistic neural tracking, we carefully controlled for neural responses related to acoustic and lexical characteristics of the speech. As pointed out by Brodbeck and Simon (2020) and Gillis et al. (2021), it is important to control for these characteristics when investigating linguistic neural tracking, as otherwise spurious linguistic neural tracking can be observed because of the high correlation between linguistic and acoustic representations. To evaluate linguistic neural tracking, we determined the added value of linguistic representations by subtracting the performance of the model containing acoustic and

lexical segmentation characteristics of the speech from the performance of the model that included the same representations together with the linguistic representations. We used four linguistic representations: phoneme surprisal, cohort entropy, word surprisal, and word frequency, which, according to Gillis et al. (2021), have an added value over and beyond acoustic representations.

All linguistic representations are one-dimensional arrays with impulses at the onsets of phonemes or words. The amplitude of an impulse represents the amount of linguistic information conveyed by the phoneme or word. To obtain the timing of the phonemes and words, we used the forced aligner of the speech alignment component of the reading tutor (Duchateau et al., 2009). Similar to linguistic representations, lexical segmentation representations are one-dimensional arrays. However, the impulses' amplitudes are one and thus independent of the amount of linguistic information.

Phoneme surprisal and cohort entropy are two metrics that summarize upcoming speech content's information value and expected information value. Phoneme surprisal is calculated as the negative logarithm of the inverse conditional probability of the phoneme given the preceding phonemes of the word. It is a measure of phoneme prediction error as it represents how surprising a phoneme is given the previously uttered phonemes (Brodbeck et al., 2018; Gwilliams and Davis, 2022). Therefore, phoneme surprisal reflects to what extent the upcoming phoneme can be predicted, which is intrinsically related to the information contributed by the phoneme. The more likely a specific phoneme is, the lower its surprisal value, and the less information is gained (and vice versa). Cohort entropy is calculated as the Shannon entropy of the active cohort of words (i.e., the cohort of words congruent with the already uttered phonemes). Cohort entropy is a measure of uncertainty: higher cohort entropy implies that the active cohort consists of a larger number of words with a similar probability of being uttered (Brodbeck et al., 2018; Gwilliams and Davis, 2022). Both measures have a similar trend: phoneme surprisal and cohort entropy tend to decrease when more phonemes of the word are pronounced because the prediction of the upcoming word can be made with a higher precision. However, if the upcoming word cannot be predicted with high precision, these metrics at phoneme level can increase.

To calculate both representations, we used a custom pronunciation dictionary that maps a word to its phoneme representation. This dictionary was created by manual and grapheme-to-phoneme conversion and contained the segmentation of 9157 words. The word probabilities were derived from the SUBTLEX-NL database (Keuleers et al., 2010). The linguistic information of the initial phoneme was not modeled in these representations. More details regarding phoneme surprisal and cohort entropy, as well as the mathematical determinations, can be found in Brodbeck et al. (2018).

The linguistic information conveyed by a word is described by word surprisal and word frequency. Similarly as phoneme surprisal, word surprisal reflects to what extent the upcoming word can be predicted based on the information contributed by that word. We used a 5 g model to determine the negative logarithm of the conditional probability of the word given the preceding words. Therefore, a word's surprisal is estimated given its four preceding words. Word frequency was derived from the same 5 g model but without including previous words, describing the word's unigram probability. This measure reflects the probability of the word, independent of its context.

These linguistic representations rely on the statistical probabilities in language. Gwilliams and Davis (2022) reviewed in detail how listeners rely on this type of internal language model that allows them to predict the upcoming word or phoneme.

Determination of neural tracking. To determine neural tracking, we used a forward modeling approach, estimating how the brain responds to specific speech characteristics. The temporal response function (TRF) describes the relationship between the presented stimulus and measured EEG. It also allows us to predict the EEG responses associated with the speech stimulus. By correlating the predicted EEG responses with the measured EEG responses, we obtain a prediction accuracy per EEG channel. This prediction accuracy is a measure of neural tracking.

We used the boosting algorithm (David et al., 2007) implemented by the Eelbrain Toolbox (Brodbeck, 2020) to estimate the TRF and obtain the prediction accuracy. The boosting algorithm estimates the TRF in a sparse and iterative way. It requires a training partition to estimate the TRF and a validation partition to evaluate the TRF. The boosting algorithm aims to find the optimal solution by evaluating the effect of incremental changes in the TRF. The incremental change, which yields the highest prediction accuracy on the validation set, is kept. Incremental changes are made until no further improvement can be found on the validation set.

We used an integration window of -100 to 600 ms; that is, the neural response is estimated ranging from 100 ms before activation of the stimulus characteristic to 600 ms after its activation. We use a broad integration window to ensure that the model captures the brain responses to the linguistic representations, which occur at longer latencies. As each speech rate condition was presented twice, we estimated the TRF on the concatenation of these two blocks per speech rate (i.e., 10 min of data). Before the TRF estimation, the data are normalized by dividing by the Euclidean norm per channel. We applied this normalization for the stimulus and EEG data individually. Then the boosting algorithm estimates the associated response behavior using a fixed step size of 0.005 . We derived the TRF and prediction accuracy per channel using a cross-validation scheme: the TRF was estimated and validated on partitions unseen during testing of the TRF to obtain the prediction accuracy. More specifically, we used 10 -fold cross-validation, implying that the data were split into 10 equally long folds, of which eight folds are used for estimating the TRF, one fold for validation, and one fold for testing. The obtained TRFs and prediction accuracies are then averaged across the different folds. The validation fold is required to determine the stopping criterion: we used an early stopping based on the ℓ_2 -norm; that is, estimation of the TRF is stopped when the Euclidean distance between the actual and predicted EEG data on the validation partition stops decreasing. The resulting TRFs are sparse. Therefore, to account for the intersubject variability and obtain a meaningful average TRF response across subjects, we smoothed the TRFs across time by convolving the estimated response with a hamming window of 50 ms in the time dimension.

To determine the acoustic tracking of the speech, we purely used acoustic representations. Therefore, we determined the prediction accuracy and TRFs based on the spectrogram and acoustic edges. Regarding the linguistic tracking of speech, we investigated the added value of these linguistic representations. To determine the added value, we subtracted the prediction accuracies of two different models. First, we estimated a baseline model consisting of acoustic and lexical segmentation representations. Second, we estimated a combined model, which included linguistic representations on top of the acoustic and lexical segmentation representations. By subtracting the prediction accuracy obtained with the baseline model from the prediction accuracy of the combined model, we can examine the added value of the linguistic representations after controlling for the acoustic and lexical segmentation representations. We did not include lexical features in the acoustic or linguistic model because the neural response to a phoneme or word onset does not reflect purely acoustic nor linguistic processing (Sanders and Neville, 2003).

We used two predetermined channel selections to investigate the effect of acoustic and linguistic tracking. The neural responses to acoustics are significantly different from those to linguistic content and therefore require a different channel selection. We used a frontal channel selection for acoustic neural tracking based on Lesenfans et al. (2019) and a central channel selection for linguistic neural tracking as reported by Gillis et al. (2021).

These channel selections were used to visualize the TRFs and to determine associated peak latency and amplitudes. To determine the peak characteristics, we set a preset time window based on the TRF averaged across subjects (Table 1). Within this time window, we normalized the TRF per channel by dividing the TRF by its ℓ_2 -norm over time to decrease across subject variability and averaged the TRF across the channel selection. Depending on a positive or negative peak, we determined the maximal or minimal amplitude and its corresponding latency to obtain the peak amplitude and latency. If the peak latency

Table 1. Time windows selected per speech representation to determine the peak characteristics

Speech representation	Time window(s) (ms)	Channel selection
Spectrogram & acoustic edges	0-90, 110-200	Frontocentral
Phoneme surprisal & cohort entropy	200-300	Central
Word surprisal & word frequency	350-450	Central

was the same as the beginning of the window, indicating the end of the previous peak, we discarded the peak from the analysis (Table 2).

Statistics

Statistical analysis was performed using MATLAB (version R2018a) and R (version 3.4.4) software. The significance level was set at $\alpha = 0.05$ unless otherwise stated.

To evaluate the subjectively rated speech understanding results, we calculated the correlation between speech rate and rated speech understanding using a Spearman rank correlation. In addition, we fitted a sigmoid function on the data to address the relation between rated speech understanding and speech rate using the `minpack.lm` package (Elzhov et al., 2016) in R. For further statistical analysis, we selected speech rate (and not subjectively rated speech understanding) as a main predictor. We opted for this because (1) it is difficult to obtain reliable speech understanding scores from a continuous story. In clinical practice, participants are asked to repeat sentences or words to calculate a percentage correct. Asking someone to repeat an entire discourse and scoring it is not feasible (MacPherson and Akeroyd, 2014; Decruy et al., 2018). (2) A subjective rating is very idiosyncratic. This results in a high inter-subject variability: some participants will rate their speech understanding more optimistically than others for the same underlying level of understanding.

To determine whether the topographies or the TRFs were significantly different from zero, we performed nonparametric permutation tests (Maris and Oostenveld, 2007). For the analysis of the acoustic TRFs, we limited the window of interest to the time region between 0 and 200 ms. As speech is more difficult to understand, the latency of the neural responses to acoustic representation increases (Mirkovic et al., 2019; Gillis et al., 2022a; Verschuere et al., 2021). These effects are most prominent in a time region of 0-200 ms (Mirkovic et al., 2019; Verschuere et al., 2020; Kraus et al., 2021), explaining the rationale to limit the time window of interest. However, no time window of interest was set to determine the significance of the linguistic TRFs. We are not aware of any studies that assess the effect of linguistic tracking when speech comprehension becomes challenging. Therefore, we chose not to specify a time window of interest when investigating the neural responses to linguistic representations. As observed in previous literature, linguistic TRFs are associated with negative responses in central areas. Therefore, we applied this test in a one-sided fashion; that is, we determined where the TRF was significantly negative.

To assess the relationship between speech rate, neural speech tracking, and speech understanding, we created a linear mixed effect (LME) model using the `LME4` package (Bates et al., 2015) in R with the following general formula:

$$\text{neuralMeasure} \sim \text{rate}(+\text{rate}^2)(+\text{understanding}) \\ + \text{random} = \text{participant}$$

where “neuralMeasure” refers to neural speech tracking, TRF amplitude or TRF latency, depending on the model being investigated and “rate” refers to the speech rate at which the speech was presented. Speech rate was also added as a quadratic effect, “rate²,” as we do not expect neural speech tracking will decrease or increase linearly indefinitely with increasing speech rate. Last, “understanding,” referring to rated speech understanding, was added to the model to investigate whether speech understanding is able to explain additional variance on top of speech rate. An additional random intercept per participant was included in the model to account for the multiple observations per

participant. “Rate²” and “understanding” are added between brackets to the general formula because these factors were only included if they benefited the model. We controlled this by calculating the Akaike Information Criterion (AIC) for the model with and without “Rate²” and “understanding.” The model with the lowest AIC was selected and its residuals plot was analyzed to assess the normality assumption of the LME residuals. Unstandardized regression coefficients (β) with 95% CIs and p value of the factors included in the model are reported in Results.

Results

Effect of speech rate on speech understanding

Figure 2 shows that, when speech rate increases, rated speech understanding decreases ($r = -0.91$, $p < 0.001$, Spearman rank correlation). To model the data, we fitted a sigmoid function between speech understanding and speech rate. The function shows a plateau until 6 syllables/s ($= \pm 1.5$ times the rate of normal speech). When the speech rate further increases, speech understanding drops. Because speech understanding and speech rate are highly correlated, selecting the one or the other for the analysis should not matter too much. We select speech rate for further analysis in function of neural speech tracking. As mentioned above, speech rate is more reliable than the subjectively rated speech understanding scores since it is objectively derived from the acoustic stimulus.

Effect of speech rate on neural processing of speech

To obtain the results in this section, we created two models: an acoustic model and a linguistic model, as explained in detail in Determination of neural tracking.

Acoustic neural tracking

First, we investigated the acoustic model containing acoustic edges and the spectrogram. Figure 3A shows how accurately the acoustic model can predict the speech signal for every electrode used. To better quantify this result, we selected the frontocentral channels based on Lesenfants et al. (2019) (channels are highlighted in red in the inset of Fig. 3B) and averaged them per subject. This resulted in one neural tracking value per speech rate per subject. Neural tracking of this frontocentral channel selection increased with increasing speech rate ($p < 0.001$, $b = 7.61 \times 10^{-3}$, CI (95%): $\pm 1.59 \times 10^{-3}$, LME; Table 3). However, as visualized on the right in Figure 3A, this increase is not monotonous, but quadratic ($p < 0.001$, $b = -3.28 \times 10^{-4}$, CI (95%): $\pm 8.41 \times 10^{-5}$, LME, Table 3). Finally, adding speech understanding as an extra predictor to the model does not improve the model ($\text{AIC}_{\text{SR}} = -623$, $\text{AIC}_{\text{SR} + \text{understanding}} = -605$).

To better understand the obtained quadratic tendency of neural tracking as a function of speech rate, we analyzed the acoustic features separately using TRFs. Figure 3B visualizes the averaged TRF of the frontocentral channels for the spectrogram (left panel) and for the acoustic edges (right panel). For both speech features, two significant positive peaks appear at ~ 70 and 150 ms (horizontal bars show the TRF parts significantly different from zero). The topographies of these peaks are shown underneath the TRFs. More detailed analysis on both peaks is done by calculating the maximum value for every participant per speech rate between 0 and 90 ms (= peak value 1) and between 110 and 200 ms (= peak value 2). We investigated the amplitude and the latency of these peak values as a function of speech rate as shown in Figure 4. With increasing speech rate, the amplitude of Peak 1 increases, saturates, and decreases again for the spectrogram

Table 2. Number of peaks (n) detected per speech representation per speech rate with $n_{\max} = 18$ (= amount of participants)

	2.6 syllables/s	3.6 syllables/s	6.2 syllables/s	9.0 syllables/s	12.9 syllables/s	16.2 syllables/s
Spectrogram: Peak 1	15	17	18	18	18	18
Spectrogram: Peak 2	14	16	18	15	13	11
Acoustic edges: Peak 1	17	17	18	18	18	18
Acoustic edges: Peak 2	18	17	16	10	8	8
Phoneme surprisal	17	18	17	16	15	18
Cohort entropy	17	16	16	15	18	17
Word surprisal	15	16	15	18	16	17
Word frequency	16	15	13	14	15	17

feature (SR: $p < 0.001$, $b = 3.95 \times 10^{-2}$, CI (95%): $\pm 8.00 \times 10^{-3}$; SR²: $p < 0.001$, $b = -1.45 \times 10^{-3}$, CI (95%): $\pm 4.20 \times 10^{-4}$; LME; Table 3) and acoustic edges (SR: $p < 0.001$, $b = 2.33 \times 10^{-2}$, CI (95%): $\pm 1.05 \times 10^{-2}$; SR²: $p < 0.001$, $b = -1.33 \times 10^{-3}$, CI (95%): $\pm 5.54 \times 10^{-4}$; LME; Table 3). In contrast, the amplitude of Peak 2 decreases with increasing speech rate (spectrogram: $p < 0.001$, $b = -6.64 \times 10^{-3}$, CI (95%): $\pm 1.77 \times 10^{-3}$; acoustic edges: $p < 0.001$, $b = -3.92 \times 10^{-3}$, CI (95%): $\pm 1.83 \times 10^{-3}$; LME; Table 3). Interestingly, the second peak for the acoustic edges even disappears when the speech rate is 9 syllables/s or higher and speech understanding drops to <80% (Fig. 3B, horizontal bars show the TRF parts significantly different from zero). For the latency analysis of Peak 2 for acoustic edges, we thus only include the latency of the peaks in the three easiest speech rate conditions, as no peaks (and latencies) can be found anymore at higher speech rates. The latency of Peak 1, for both speech features, increases with increasing speech rate (spectrogram: $p < 0.001$, $b = 1.02 \times 10^{-3}$, CI (95%): $\pm 5.35 \times 10^{-4}$; acoustic edges: $p < 0.001$, $b = 1.05 \times 10^{-3}$, CI (95%): $\pm 4.58 \times 10^{-4}$; LME; Table 3), while the latency of Peak 2 shows no significant relation with speech rate (spectrogram: $p = 0.80$, $b = -9.81 \times 10^{-5}$, CI (95%): $\pm 7.55 \times 10^{-4}$; acoustic edges: $p = 0.11$, $b = 2.37 \times 10^{-3}$, CI (95%): $\pm 2.84 \times 10^{-3}$; LME; Table 3).

Linguistic neural tracking

Next to acoustic neural tracking, we also investigated the effect of speech rate on linguistic neural tracking (for more details, see Determination of neural tracking). Figure 5A (left) shows how accurately the linguistic model can predict the speech signal over subjects per speech rate per channel. The channels in the cluster, which drives the topography from 0, are annotated with gray markers. The higher the speech rate, the fewer channels have significant neural tracking. To quantify this, we averaged the prediction accuracy over a central channel selection based on Gillis et al. (2021) (channels are highlighted in red in the inset of Fig. 5B), resulting in one neural tracking value per speech rate per subject. As shown in Figure 5A (right), neural tracking significantly drops monotonically with increasing speech rate ($p = 0.008$, $b = -4.59 \times 10^{-5}$, CI (95%): $\pm 3.31 \times 10^{-5}$, LME; Table 4). Interestingly, this is the opposite trend from the acoustic model in Acoustic neural tracking. Finally, adding speech understanding as a predictor does not improve the linguistic model ($AIC_{SR} = -1175$, $AIC_{SR + understanding} = -1151$).

To thoroughly investigate the neural responses to the linguistic features, we examined the TRFs of the central channel selection. Figure 5B visualizes the averaged normalized TRF in the central channel selection for the different linguistic features per speech rate. Significant responses can be found in the lower speech rates when speech can still be understood for all features.

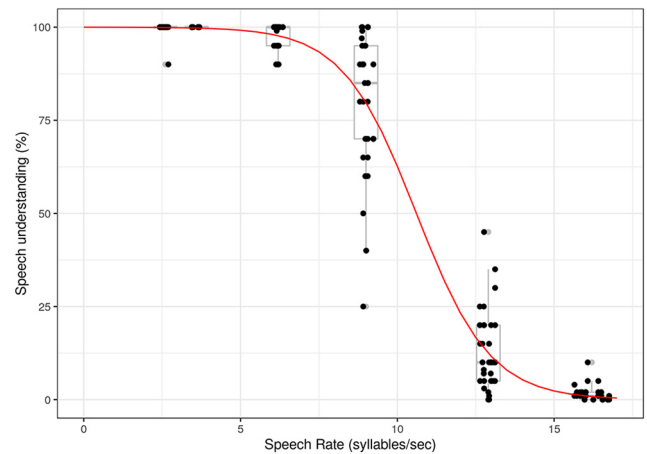
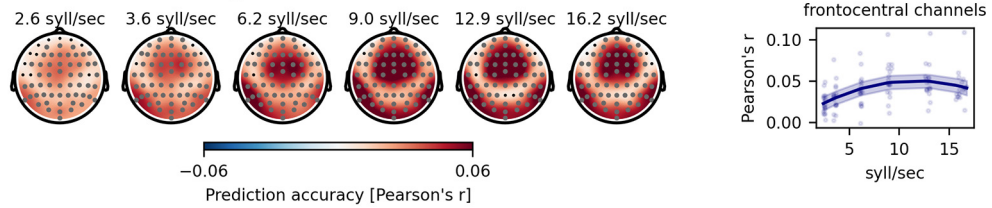


Figure 2. Rated speech understanding in function of speech rate. Dots represent speech understanding per participant for every participant-specific speech rate. Boxplots represent the participants' results for the averaged speech rates. Red line indicates the sigmoid function fitted on the data over participant.

In the higher speech rates, where speech understanding is worse or absent, the linguistic neural response also disappears (Fig. 5B, horizontal bars show the TRF parts significantly different from zero). The topographies of these responses are shown in Figure 5B underneath the TRFs. Interestingly, the topographies switch from central negativity when speech is understood to frontal negativity when speech understanding is worse or absent. To investigate whether the amplitude or latency of these peaks is related to speech rate, we calculated the minimum value for every participant within the corresponding time region, defined in Table 1 (= peak value). For all linguistic features, the peak amplitude shrinks significantly with increasing speech rate as shown in Figure 6 (phoneme surprisal: $p < 0.001$, $b = 6.81 \times 10^{-3}$, CI (95%): $\pm 2.82 \times 10^{-3}$; cohort entropy: $p = 0.008$, $b = 4.21 \times 10^{-3}$, CI (95%): $\pm 3.07 \times 10^{-3}$; word surprisal: $p < 0.001$, $b = 6.46 \times 10^{-3}$, CI (95%): $\pm 3.09 \times 10^{-3}$; word frequency: $p < 0.001$, $b = 6.44 \times 10^{-3}$, CI (95%): $\pm 3.00 \times 10^{-3}$; LME; Table 4). In other words, when speech becomes faster and more difficult to understand, the peak amplitude of the linguistic features decreases until it finally disappears when the speech rate is 9.0 or 12.9 syllables/s, or higher, and speech understanding is dropping below 90%. Similar to the analysis of the second peak for acoustic features, we only include the latency of significant peaks. Phoneme surprisal shows a significant increase of peak latency with increasing speech rate ($p = 0.015$, $b = 5.47 \times 10^{-3}$, CI (95%): $\pm 4.20 \times 10^{-3}$; LME; Table 4). Cohort entropy, word surprisal and word frequency, on the other hand, reveal no significant effect of speech rate on peak latency (cohort entropy: $p = 0.18$, $b = 3.59 \times 10^{-3}$, CI (95%): $\pm 5.11 \times 10^{-3}$; word surprisal: $p = 0.23$, $b = -3.33 \times 10^{-3}$, CI

A Acoustic neural tracking



B Temporal Response Function

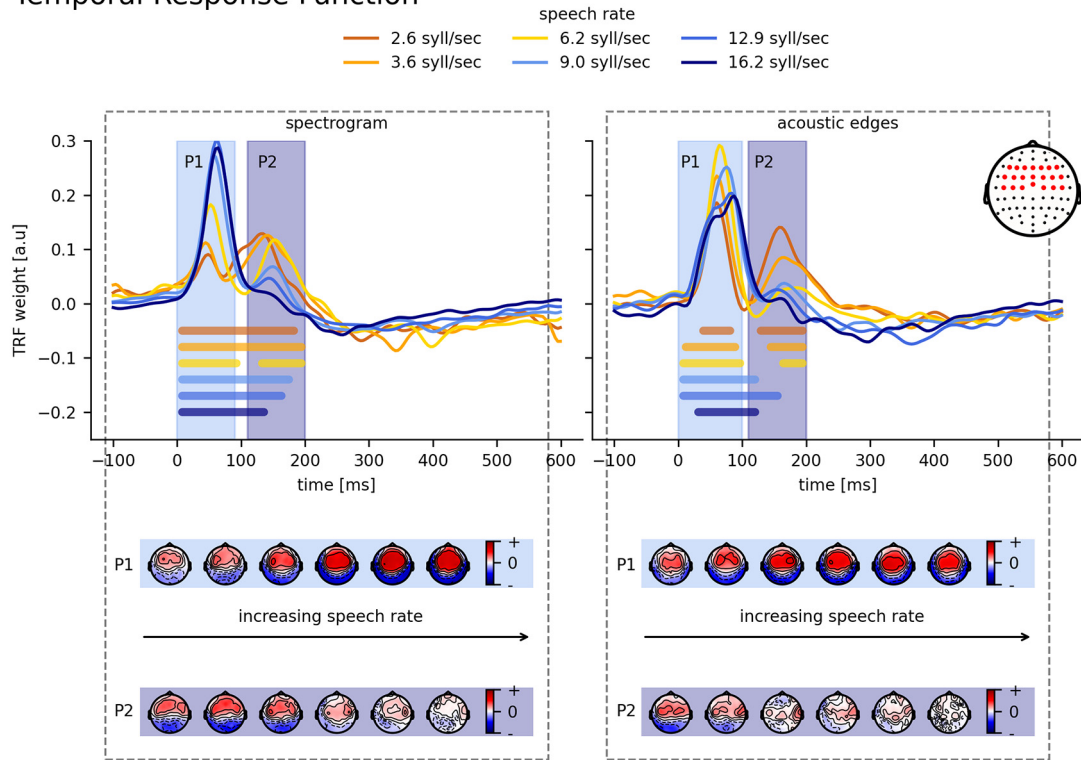


Figure 3. Effect of speech rate on acoustic tracking. **A**, Visualization of the average prediction accuracy across participants for each speech rate. Annotated gray channels represent the cluster that drives the significant difference from 0. Right, How acoustic tracking, averaged across frontocentral channels, changes according to the speech rate. **B**, Normalized TRFs of the spectrogram and acoustic edges. Bold horizontal lines indicate where the TRFs are significantly different from 0 (the same color as the TRF of the considered speech rate). The topographies below represent the associated peak topographies in the TRF for the first (P1) and second peak (P2).

Table 3. LME model of prediction accuracies and amplitude and latency of the TRF peaks in function of speech rate for the acoustic model^a

			Rate			Rate ²		
			β	CI (95%)	<i>p</i>	β	CI (95%)	<i>p</i>
Prediction accuracy			7.61×10^{-3}	$\pm 1.59 \times 10^{-3}$	<0.001	-3.28×10^{-4}	$\pm 8.41 \times 10^{-5}$	<0.001
Spectrogram	p1	TRF amplitude	3.95×10^{-2}	$\pm 8.00 \times 10^{-3}$	<0.001	-1.45×10^{-3}	$\pm 4.20 \times 10^{-4}$	<0.001
		TRF latency	1.02×10^{-3}	$\pm 5.35 \times 10^{-4}$	<0.001	Does not improve AIC		
	p2	TRF amplitude	-6.64×10^{-3}	$\pm 1.77 \times 10^{-3}$	<0.001	Does not improve AIC		
		TRF latency	-9.81×10^{-5}	$\pm 7.55 \times 10^{-4}$	NS	Does not improve AIC		
Acoustic edges	p1	TRF amplitude	2.33×10^{-2}	$\pm 1.05 \times 10^{-2}$	<0.001	-1.33×10^{-3}	$\pm 5.54 \times 10^{-4}$	<0.001
		TRF latency	1.05×10^{-3}	$\pm 4.58 \times 10^{-4}$	<0.001	Does not improve AIC		
	p2	TRF amplitude	-3.92×10^{-3}	$\pm 1.83 \times 10^{-3}$	<0.001	Does not improve AIC		
		TRF latency	2.37×10^{-3}	$\pm 2.84 \times 10^{-3}$	NS	Does not improve AIC		

^aEvery line indicates a different model. p1, first peak; p2, second peak.

(95%): $\pm 5.38 \times 10^{-3}$; word frequency: $p = 0.95$, $b = -1.09 \times 10^{-4}$, CI (95%): $\pm 3.06 \times 10^{-3}$; LME; Table 4).

found that neural tracking of the acoustic features increased, while neural tracking of the linguistic features decreased.

Discussion

We aimed to investigate whether neural speech processing can capture the effect of gradually decreasing speech understanding by manipulating the speech rate. With increasing speech rate, we

Effect of speech rate on acoustic neural processing of speech

We hypothesized that neural tracking of acoustic features would increase with increasing speech rate because more acoustic information is present. We found an increase of acoustic neural

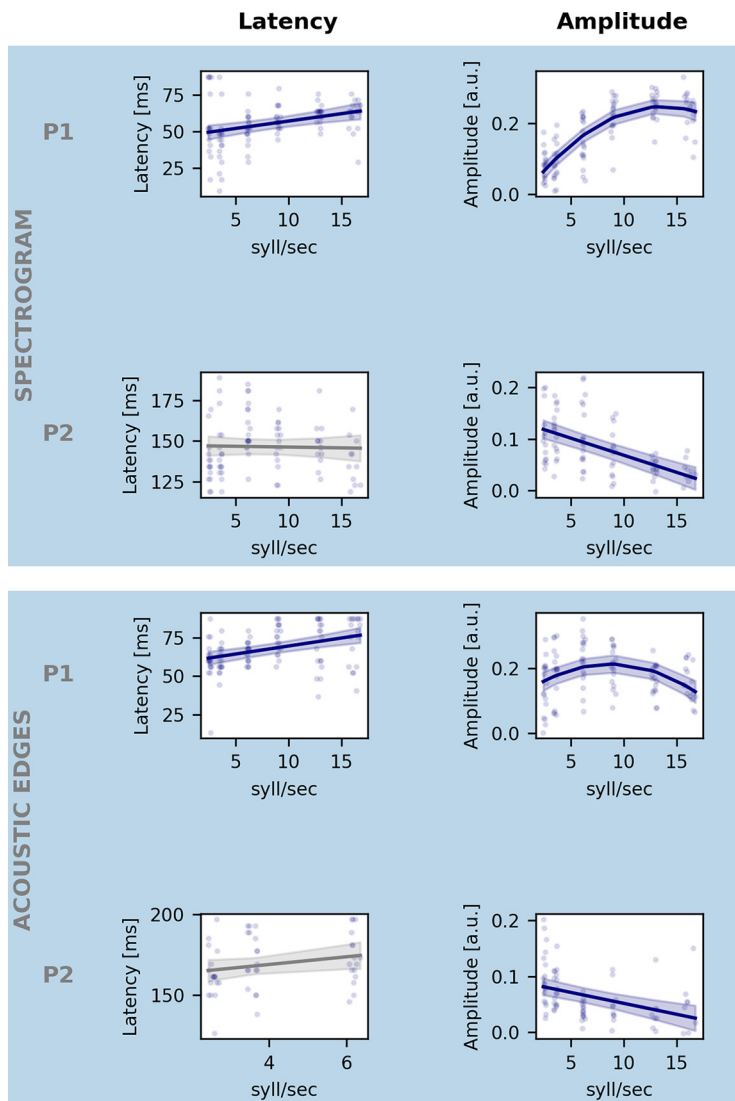


Figure 4. Effect of speech rate on the amplitude and latency of acoustic representation. Each row represents the effect of speech rate on latency (left) and the amplitude (right) of the neural response to spectrogram (top row) and acoustic edges (bottom row) for, respectively, the first (P1) and second (P2) identified peak as indicated in Figure 3. Blue line indicates the model's prediction for each speech rate. Shaded area represents the CI of the model's prediction. Gray represents the nonsignificant models. We only include the latency of significant peaks for the latency analysis.

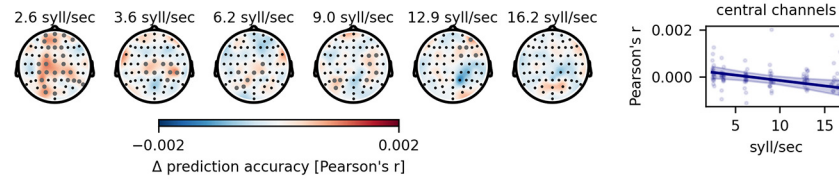
tracking with increasing speech rate and thus decreasing speech understanding, confirming our hypothesis. When speech becomes faster, the model is better at predicting acoustic speech features.

This increase of acoustic neural tracking with increasing speech rate, and thus decreasing speech understanding, seems discrepant with previous research trying to link acoustic neural tracking to speech understanding using, for example, speech-in-noise paradigms (Ding and Simon, 2013; Vanthornhout et al., 2018; Etard and Reichenbach, 2019; Iotzov and Parra, 2019; Verschuere et al., 2020). The experimental paradigm could explain this discrepancy. Previous studies used, for example, noise to manipulate speech understanding. In those cases, decreased neural tracking was accompanied by a decrease in speech understanding and an acoustically degraded speech signal. Because speech understanding and signal-to-noise ratio are highly correlated, it is challenging to unravel to what extent the decreased neural tracking is driven by decreased speech understanding or the signal-to-noise ratio used to vary speech understanding. In

this study, we manipulated speech understanding by speeding up the speech signal and preserving its signal-to-noise ratio, in contrast to the speech in noise studies. We hypothesize that the brain mainly responds to acoustic boundaries (i.e., onsets of sounds), which are more prominent in the faster speech presented in this study, explaining the increasing tendency. In addition, because the speech is sped up, the duration of the silences in between words or sentences inherently decreases, which increases the amount of speech data allowing the model to improve its estimate of the TRF and obtain higher prediction accuracies. However, this increase of acoustic neural tracking and speech rate is not linear: it saturates, and decreases again at very high speech rates. A first hypothesis is related to the motor cortex. When participants listen to the speech, they tend to mimic the speech in their brain, activating motor areas (Casas et al., 2021). However, most speakers cannot produce speech as fast as 16.2 syllables/s. Hence, the corresponding mouth movements are unnatural, which implies that the listener cannot mimic the speech in their brain anymore, decreasing the related responses in the motor areas and thus brain responses to the acoustic speech features in general. A different hypothesis is related to the stimulus characteristics. When increasing the speech rate, the spectrogram and acoustic edges contain more and more peaks. Possibly, at very high speech rates, these peaks occur so fast after each other, making it difficult for the brain to perceive them still (see Fig. 1), similar to speech in noise studies where the amount of noise can mask acoustic boundaries. Therefore, it is difficult to attribute this decrease in acoustic neural tracking at very high speech rates (or low signal-to-noise ratios): a decrease in speech understanding or a decrease in neural detection of acoustic boundaries, or a combination of both?

To better understand the observed quadratic tendency of acoustic neural tracking with increasing speech rate, we investigated the TRFs of the speech features separately. Two significant peaks with opposite behavior could be observed for both acoustic features. The first peak is the largest. With increasing speech rate, the amplitude of Peak 1 increases, saturates, and decreases again at very high speech rates. This is similar to the previously discussed neural acoustic tracking results. On the other hand, the second peak amplitude decreases with increasing speech rate. This discrepancy is intriguing as it suggests that both peaks have different underlying brain processes as confirmed by literature (Picton, 2011; Brodbeck and Simon, 2020). Peak 1 occurs relatively fast, at ~50 ms, and is probably mainly related to the acoustics of the incoming speech and thus benefits from an increased speech rate. Peak 2, on the other hand, occurs somewhat later, at ~150 ms, and could, in addition to the acoustics, be influenced by top-down processing related to speech understanding and attention (Ding and Simon, 2012; Vanthornhout et al., 2019). In addition to the amplitude, we also investigated the latencies. The latency of the first peak increases with increasing speech rate. Increased latencies are often

A Linguistic neural tracking



B Temporal Response Function

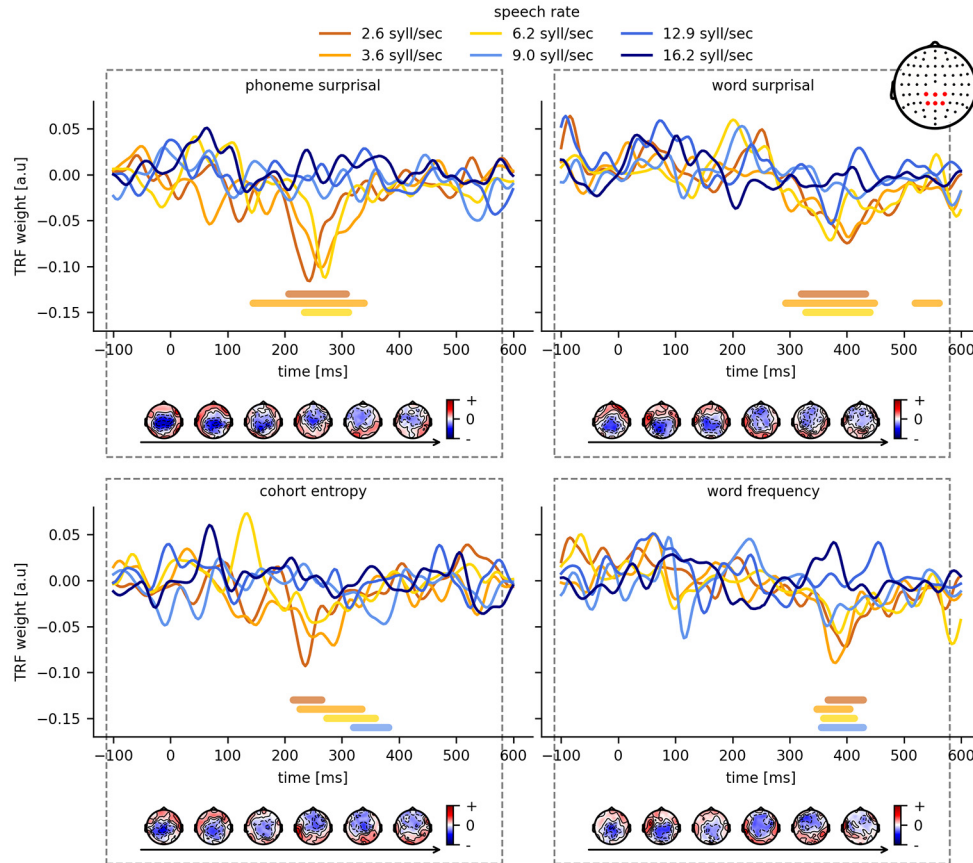


Figure 5. Effect of speech rate on linguistic tracking. **A**, Visualization of the added value of linguistic representations across participants for each speech rate. Annotated gray channels represent the cluster that drives the significant difference from 0. Right, How linguistic tracking, averaged across central channels, changes according to the speech rate. **B**, The associated normalized TRFs for the linguistic representations. Bold horizontal lines indicate where the TRFs are significantly different from 0 (the same color as the TRF of the considered speech rate). The topographies below represent the associated peak topographies to the TRF. Horizontal arrow underneath the topographies indicates the increasing speech rate.

Table 4. LME model of prediction accuracy and amplitude and latency of the TRF peaks in function of speech rate for the linguistic model^a

		Rate			Rate ²		
		β	CI (95%)	<i>p</i>	β	CI (95%)	<i>p</i>
Prediction accuracy		-4.59×10^{-5}	$\pm 3.31 \times 10^{-5}$	0.0081	Does not improve AIC		
Phoneme surprisal	TRF amplitude	6.81×10^{-3}	$\pm 2.82 \times 10^{-3}$	<0.001	Does not improve AIC		
	TRF latency	5.47×10^{-3}	$\pm 4.20 \times 10^{-3}$	0.015	Does not improve AIC		
Cohort entropy	TRF amplitude	4.21×10^{-3}	$\pm 3.07 \times 10^{-3}$	0.008	Does not improve AIC		
	TRF latency	3.59×10^{-3}	$\pm 5.11 \times 10^{-3}$	NS	Does not improve AIC		
Word surprisal	TRF amplitude	6.46×10^{-3}	$\pm 3.09 \times 10^{-3}$	<0.001	Does not improve AIC		
	TRF latency	-3.33×10^{-3}	$\pm 5.38 \times 10^{-3}$	NS	Does not improve AIC		
Word frequency	TRF amplitude	6.44×10^{-3}	$\pm 3.00 \times 10^{-3}$	<0.001	Does not improve AIC		
	TRF latency	-1.09×10^{-4}	$\pm 3.06 \times 10^{-3}$	NS	Does not improve AIC		

^aEvery line indicates a different model.

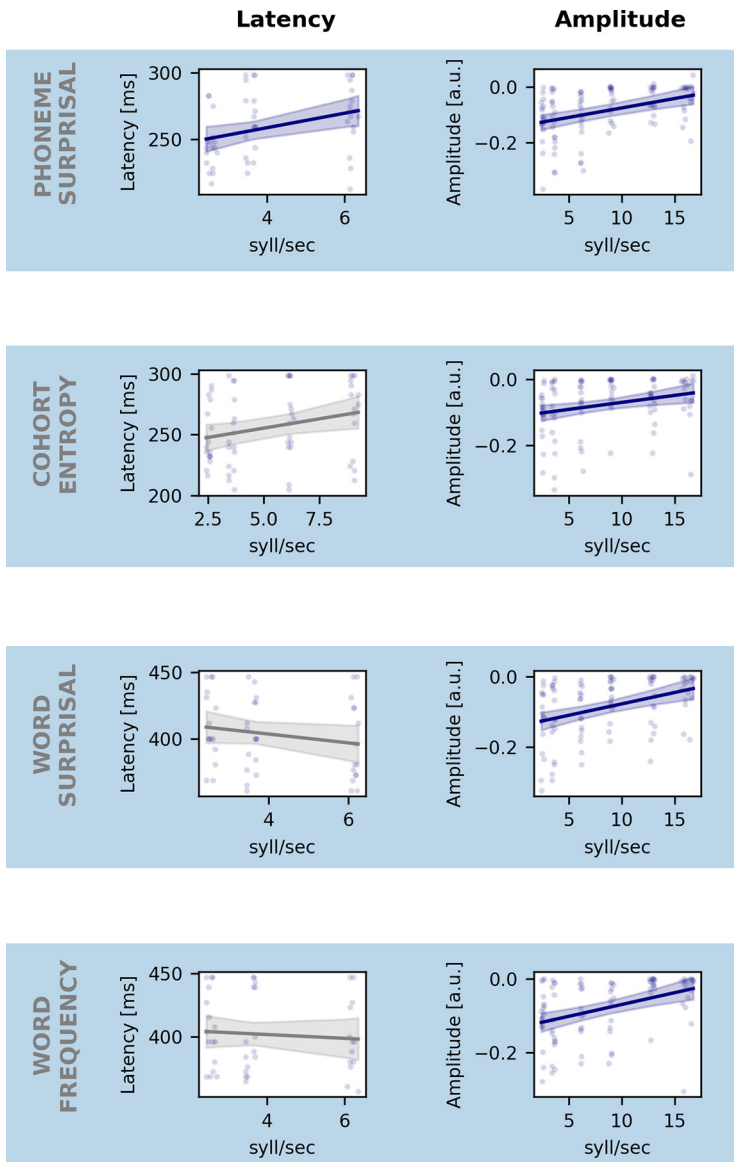


Figure 6. Effect of speech rate on the amplitude and latency of linguistic representation. Each row represents the effect of speech rate on latency (left) and the amplitude (right) of the neural response to phoneme surprisal (top row), cohort entropy (second row), word surprisal (third row), and word frequency (bottom row). Blue line indicates the model's prediction for each speech rate. Shaded area represents the CI of the model's prediction. Gray represents the nonsignificant models. We only include the latency of significant peaks for the latency analysis.

observed in more complex conditions with a higher task demand, like, for example, lower stimulus intensity, vocoded speech, or speech in noise (Mirkovic et al., 2019; Kraus et al., 2021; Verschueren et al., 2021). The latency of the neural responses can also be related to neural processing efficiency (Bidelman et al., 2019; Gillis et al., 2022a). In more detail, a larger latency indicates that more processing time is required to process the same speech characteristics, showing reduced neural processing efficiency. More words and phonemes need to be processed as the speech rate increases, resulting in a more challenging condition to process the incoming speech.

Effect of speech rate on linguistic neural processing of speech

When speech becomes faster, speech understanding drops. Interestingly, this same decrease can be observed in linguistic neural tracking, in contrast to acoustic neural tracking (see Effect of speech rate on acoustic neural processing of speech). It is

interesting to note that the pattern between speech rate and acoustic and linguistic neural tracking differs. Linguistic neural tracking seems linearly related to speech rate, while acoustic neural tracking shows nonlinearity: it increases, stagnates, and eventually slightly decreases for the highest speech rates. To the best of our knowledge, this is the first study that evaluates linguistic neural tracking when manipulating the level of speech understanding as a gradual effect. The studies of Brodbeck et al. (2018) and Broderick et al. (2018) using a two-talker paradigm are most comparable. They compared two conditions (i.e., intelligible and attended speech vs unintelligible and ignored speech), but not the spectrum in between. Nevertheless, their findings converge with our results and support our hypothesis of linguistic neural tracking as a neural marker of speech understanding. When the speech is not understood or ignored, the brain does not track the linguistic aspects of the speech, while for intelligible speech linguistic tracking is present.

To better understand the observed decrease of linguistic neural tracking with increasing speech rate, we investigated the TRFs of the speech representations separately. We observed a characteristic negative peak for each linguistic representation as observed in previous literature (see, e.g., Brodbeck et al., 2018; Weissbart et al., 2020; Gillis et al., 2021). Although the responses are very small, they are significantly different from zero and the TRF patterns converge with the patterns observed in previous literature. For the phoneme-related features (phoneme surprisal and cohort entropy), this peak occurs at ~250 ms. For the word-related features (word surprisal and word frequency), this peak occurs somewhat later, ~350 ms. The difference in timescale between both feature groups could be linked to the different speech processing stages (phonemes vs words) they represent (Gillis et al., 2022b). Regarding the topographies of these peaks, the understandable speech conditions are associated with a typical topography, similar to the classical N400 responses characterized by central negative channels. As speech becomes less understandable (i.e., the speech rate increases), the associated topography disappears.

For all speech features, the amplitude of this negative peak decreases with increasing speech rate until it disappears at speech rates as high as 9.0–12.9 syllables/s. Gillis et al. (2021) already showed that these linguistic representations have an added value above and beyond acoustic and lexical representations. However, the authors did not compare intelligible to unintelligible speech. Here, we elegantly showed that, as the speech becomes less understandable but remains audible and acoustically intact (in contrast to speech in noise studies or vocoder studies), the characteristic negative peak decreases and finally disappears. Together, our results suggest that these characteristic negative peaks to linguistic representations could be neural correlates of speech understanding.

Caveats of this study

As mentioned in the Discussion, it is difficult to fully attribute a decrease in acoustic neural tracking at very high speech rates to a decrease in speech understanding (top-down) or a decrease in neural detection of acoustic boundaries (bottom-up). This is an important caveat in studies manipulating speech understanding by changing the acoustics of the speech signal. We cannot be sure whether the observed decrease is driven by speech rate or speech understanding. To enable a clean separation of acoustics and speech understanding, the speech acoustics should remain the same over conditions, such as in the paradigm used by Di Liberto et al. (2018). They use a priming paradigm where vocoded speech is presented before and after priming (i.e., learning the speech's meaning by presenting the nonvocoded sentence before the vocoded sentence). This is an elegant way to disentangle acoustic from linguistic processing but only enables an on-off (understanding vs no understanding) analysis. In addition, priming is only feasible with a small number of words or sentences. Using a continuous story would make this task too complex for the participants.

The second caveat in this study is that, per condition, neural speech tracking is trained on the same amount of data in terms of time but not feature instances. For example, 10 min of slow speech contains fewer feature instances than 10 min of fast speech. In this study, we elected to keep the amount of EEG data constant because we know a model is affected by the amount of training data. Therefore, it is hard to conclude to what extent the increased acoustic neural tracking at higher speech rates is influenced by the increased number of feature instances in the model. Ideally, we would be able to compare both models. This is unfortunately not convenient in this study because 10 min of slow speech corresponds to only 94 s of speech in the fastest condition, which is not enough to train a model. However, as suggested by the reviewer, we did redo the analysis for the four lowest speech rates reducing the amount of EEG until each speech rate condition contained the same amount of feature instances. The effect of speech rate on acoustic tracking was similar for both analyses: when speech rate increases, acoustic neural tracking increases. The only difference is that this increase is faster when keeping the amount of data constant instead of preserving the amount of feature instances.

In conclusion, using a speech rate paradigm, we map how the level of speech understanding affects acoustic and linguistic neural speech processing. When speech rate increases, acoustic neural tracking increases, although speech understanding drops. However, the amplitude of the later acoustic neural response decreases with increasing speech rate, suggesting influence of top-down processing related to speech understanding and attention. In contrast, linguistic neural tracking decreases with increasing speech rate and even disappears when speech is no longer understood. Together, this suggests that linguistic neural tracking could possibly be a more direct predictor of speech understanding compared with acoustic neural tracking.

References

- Ahissar E, Nagarajan S, Ahissar M, Protopapas A, Mahncke H, Merzenich MM (2001) Speech comprehension is correlated with temporal response patterns recorded from auditory cortex. *Proc Natl Acad Sci USA* 98:13367–13372.
- Baayen RH, Piepenbrock R, Gulikers L (1996) The celex lexical database (cd-rom).
- Bates D, Mächler M, Bolker B, Walker S (2015) Fitting linear mixed-effects models using lme4. *J Stat Softw* 67:1–48.
- Bidelman GM, Price CN, Shen D, Arnott SR, Alain C (2019) Afferent-efferent connectivity between auditory brainstem and cortex accounts for poorer speech-in-noise comprehension in older adults. *Hear Res* 382:107795.
- Boursma P, Weenink D (2018) Praat: doing phonetics by computer [computer program], version 6.0.37. Available at <http://www.praat.org/>.
- Brodbeck C (2020) Eelbrain 0.34. Available at <http://doi.org/10.5281/zenodo.3923991>.
- Brodbeck C, Hong LE, Simon JZ (2018) Rapid transformation from auditory to linguistic representations of continuous speech. *Curr Biol* 28:3976–3983.e5.
- Brodbeck C, Simon JZ (2020) Continuous speech processing. *Curr Opin Physiol* 18:25–31.
- Broderick MP, Anderson AJ, Liberto GM, Crosse MJ, Lalor EC (2018) Electrophysiological correlates of semantic dissimilarity reflect the comprehension of natural, narrative speech. *Curr Biol* 28:803–809.e3.
- Casas AS, Lajnef T, Pascarella A, Guiraud-Vinatea H, Laaksonen H, Bayle D, Jerbi K, Boulenger V (2021) Neural oscillations track natural but not artificial fast speech: novel insights from speech-brain coupling using MEG. *Neuroimage* 244:118577.
- David SV, Mesgarani N, Shamma SA (2007) Estimating sparse spectro-temporal receptive fields with natural stimuli. *Network* 18:191–212.
- Decruy L, Das N, Verschuere E, Francart T (2018) The self-assessed Békésy procedure: validation of a method to measure intelligibility of connected discourse. *Trends Hear* 22:2331216518802702.
- Di Liberto GM, Lalor EC, Millman RE (2018) Causal cortical dynamics of a predictive enhancement of speech intelligibility. *Neuroimage* 166:247–258.
- Ding N, Simon JZ (2012) Emergence of neural encoding of auditory objects while listening to competing speakers. *Proc Natl Acad Sci USA* 109:11854–11859.
- Ding N, Simon JZ (2013) Adaptive temporal encoding leads to a background-insensitive cortical representation of speech. *J Neurosci* 33:5728–5735.
- Duchateau J, Kong YO, Cleuren L, Latacz L, Roelens J, Samir A, Demuynck K, Ghesquière P, Verhelst W, Van Hamme H (2009) Developing a reading tutor: design and evaluation of dedicated speech recognition and synthesis modules. *Speech Commun* 51:985–994.
- Elzhov TV, Mullen KM, Spiess AN, Bolker B (2016) minpack.lm: R interface to the Levenberg-Marquardt nonlinear least-squares algorithm found in MINPACK, plus support for bounds. R package version 1.2-1.
- Etard O, Reichenbach T (2019) Neural speech tracking in the theta and in the delta frequency band differentially encode clarity and comprehension of speech in noise. *J Neurosci* 39:5750–5759.
- Francart T, van Wieringen A, Wouters J (2008) APEX 3: a multi-purpose test platform for auditory psychophysical experiments. *J Neurosci Methods* 172:283–293.
- Gillis M, Vanthornhout J, Simon JZ, Francart T, Brodbeck C (2021) Neural markers of speech comprehension: measuring EEG tracking of linguistic speech representations, controlling the speech acoustics. *J Neurosci* 41:10316–10329.
- Gillis M, Decruy L, Vanthornhout J, Francart T (2022a) Hearing loss is associated with delayed neural responses to continuous speech. *European Journal of Neuroscience* 55:1671–1690.
- Gillis M, Van Canneyt J, Francart T, Vanthornhout J (2022b) Neural tracking as a diagnostic tool to assess the auditory pathway. *bioRxiv* 470129. doi: 10.1101/2021.11.26.470129.
- Gwilliams L, Davis MH (2022) Extracting language content from speech sounds: the information theoretic approach. In: *Speech perception* (Holt LL, Peelle JE, Coffin AB, Popper AN, Fay RR, eds), Vol 74, pp 113–139. Springer, Cham.
- Heeris J (2014) Gammatone filterbank toolkit 1.0. Available at <https://github.com/detly/gammatone>.
- Hertrich I, Dietrich S, Trouvain J, Moos A, Ackermann H (2012) Magnetic brain activity phase-locked to the envelope, the syllable onsets, and the fundamental frequency of a perceived speech signal. *Psychophysiology* 49:322–334.
- Iotzov I, Parra L (2019) EEG can predict speech intelligibility. *J Neural Eng* 16:036008.
- Keuleers E, Brysbaert M, New B (2010) Subtlex-nl: a new measure for Dutch word frequency based on film subtitles. *Behav Res Methods* 42:643–650.

- Koskinen M, Kurimo M, Gross J, Hyv A, Hari R (2020) Neuroimage brain activity reflects the predictability of word sequences in listened continuous speech. *Neuroimage* 219:116936.
- Kraus F, Tune S, Ruhe A, Obleser J, Woestmann M (2021) Unilateral acoustic degradation delays attentional separation of competing speech. *Trends in Hearing* 25:23312165211013242.
- Lesenfants D, Vanthornhout J, Verschuereen E, Decruy L, Francart T (2019) Predicting individual speech intelligibility from the neural tracking of acoustic- and phonetic-level speech representations. *Hear Res* 380:1–9.
- MacPherson A, Akeroyd MA (2014) A method for measuring the intelligibility of uninterrupted, continuous speech. *J Acoust Soc Am* 135:1027–1030.
- Maris E, Oostenveld R (2007) Nonparametric statistical testing of EEG- and MEG-data. *J Neurosci Methods* 164:177–190.
- Mirkovic B, Debener S, Schmidt J, Jaeger M, Neher T (2019) Effects of directional sound processing and listener's motivation on EEG responses to continuous noisy speech: do normal-hearing and aided hearing-impaired listeners differ? *Hear Res* 377:260–270.
- Müller JA, Wendt D, Kollmeier B, Debener S, Brand T, Hunter CR (2019) Effect of speech rate on neural tracking of speech. *Front Psychol* 10:1–15.
- Nourski KV, Reale RA, Oya H, Kawasaki H, Kovach CK, Chen H, Iii MA, Brugge JF (2009) Temporal envelope of time-compressed speech represented in the human auditory cortex. *J Neurosci* 29:15564–15574.
- Picton TW (2011) Human auditory evoked potentials. San Diego: Plural.
- Sanders LD, Neville HJ (2003) An ERP study of continuous speech processing: I. Segmentation, semantics, and syntax in native speakers. *Brain Res Cogn Brain Res* 15:228–240.
- Shannon RV, Zeng FG, Kamath V, Wygonski J, Ekelid M (1995) Speech recognition with primarily temporal cues. *Science* 270:303–304.
- Slaney M (1998) Auditory toolbox: Interval Research Corporation. Tech Rep 10.
- Somers B, Francart T, Bertrand A (2018) A generic EEG artifact removal algorithm based on the multi-channel Wiener filter. *J Neural Eng* 15:036007.
- Vanthornhout J, Decruy L, Wouters J, Simon JZ, Francart T (2018) Speech intelligibility predicted from neural entrainment of the speech envelope. *J Assoc Res Otolaryngol* 19:181–191.
- Vanthornhout J, Decruy L, Francart T (2019) Effect of task and attention on neural tracking of speech. *Front Neurosci* 13:977.
- Verschuereen E, Vanthornhout J, Francart T (2020) The effect of stimulus choice on an EEG-based objective measure of speech intelligibility. *Ear Hear* 41:1586–1597.
- Verschuereen E, Vanthornhout J, Francart T (2021) The effect of stimulus intensity on neural envelope tracking. *Hear Res* 403:108175.
- Weissbart H, Kandylaki KD, Reichenbach T (2020) Cortical tracking of surprisal during continuous speech comprehension. *J Cogn Neurosci* 32:155–166.