

Linguistic modulation of the neural encoding of phonemes

Seung-Goo Kim ^{1,2,*}, Federico De Martino³, Tobias Overath^{1,4,5}

¹Department of Psychology and Neuroscience, Duke University, 308 Research Dr, Durham, NC 27708, United States

²Research Group Neurocognition of Music and Language, Max Planck Institute for Empirical Aesthetics, Grüneburgweg 14, Frankfurt am Main 60322, Germany

³Faculty of Psychology and Neuroscience, University of Maastricht, Universiteitssingel 40, 6229 ER Maastricht, Netherlands

⁴Duke Institute for Brain Sciences, Duke University, 308 Research Dr, Durham, NC 27708, United States

⁵Center for Cognitive Neuroscience, Duke University, 308 Research Dr, Durham, NC 27708, United States

*Corresponding authors: Seung-Goo Kim, Max Planck Institute for Empirical Aesthetics, Grüneburgweg 14, Frankfurt am Main 60322, Germany.

Email: dr.seunggoo.kim@gmail.com and Tobias Overath, Levine Science Research Center, Duke University Box 90999, Durham, NC 27708, United States.

Email: t.overath@duke.edu

Speech comprehension entails the neural mapping of the acoustic speech signal onto learned linguistic units. This acousto-linguistic transformation is bi-directional, whereby higher-level linguistic processes (e.g. semantics) modulate the acoustic analysis of individual linguistic units. Here, we investigated the cortical topography and linguistic modulation of the most fundamental linguistic unit, the phoneme. We presented natural speech and “phoneme quilts” (pseudo-randomly shuffled phonemes) in either a familiar (English) or unfamiliar (Korean) language to native English speakers while recording functional magnetic resonance imaging. This allowed us to dissociate the contribution of acoustic vs. linguistic processes toward phoneme analysis. We show that (i) the acoustic analysis of phonemes is modulated by linguistic analysis and (ii) that for this modulation, both of acoustic and phonetic information need to be incorporated. These results suggest that the linguistic modulation of cortical sensitivity to phoneme classes minimizes prediction error during natural speech perception, thereby aiding speech comprehension in challenging listening situations.

Key words: acoustic analysis; functional MRI; human auditory cortex; linguistic analysis; speech perception.

Introduction

Speech comprehension relies on the neural mapping of the acoustic speech signal onto linguistic categories (Hickok and Poeppel 2007; Poeppel et al. 2008; Kleinschmidt and Jaeger 2015). As such, the acoustic speech waveform that reaches our ears is converted into a neural code in the inner ear, which is then processed along the ascending auditory system and subsequently matched to learned linguistic categories (Hickok and Poeppel 2007; Friederici 2011). Importantly, the percept of phonemes requires phonemic knowledge in a given language as the variance in acoustic realization of speech sounds across utterances and speakers is too enormous to be directly handled in the acoustic domain (Liberman et al. 1967). However, while this acousto-linguistic transformation is the basis for successful speech comprehension, many aspects of it still remain unknown. For example, previous studies, using synthesized monosyllabic stimuli, have localized neural correlates of different categorization of identical speech sounds in the inferior frontal cortex and premotor cortical regions (Hasson et al. 2007; Kilian-Hütten et al. 2011; Lee et al. 2012; Preisig et al. 2022). However, it remains unclear how this knowledge translates to natural language processing. Here, we ask (i) whether the acousto-linguistic transformation is malleable to top-down linguistic information and (ii) whether we can dissociate the contributions of acoustic and linguistic processing toward this transformation.

The phoneme is the smallest perceptual unit capable of determining the meaning of a word (e.g. the words pin and chin differ

only with respect to their initial phonemes) (Stevens 2000). Of the upward of 100 phonemes in use world-wide, ~44 phonemes make up the English language and these are categorized primarily based on articulatory features into four main classes: vowels, nasals and sonorants, plosives, fricatives, and affricates (Ladefoged 2001; Ladefoged and Johnstone 2015). Each phoneme class has characteristic acoustic features; for example, while vowel sounds display a sustained period of harmonicity, plosives are characterized by a brief period of silence followed by a short broadband noise burst. Individual phonemes and the phoneme classes to which they belong have distinct temporal neural correlates: each phoneme class has a unique time-locked neural response characteristic, or phoneme-related potential (PRP; Khalighinejad et al. 2017; Overath and Lee 2017). The phoneme-class-specific PRPs likely reflect the neural analysis of their acoustic characteristics (e.g. timing of energy onset, harmonicity, etc.) in functionally separate parts of auditory cortex.

In natural speech, phonemes do not occur in isolation, but instead form sequences to create syllables and words. The order in which phonemes can occur is governed by phonotactics, and is unique to each language (Chomsky and Halle 1965). Apart from learning to recognize the language-specific phonemes themselves (Cheour et al. 1998), phonotactics is one of the first sets of rules infants need to learn during language acquisition (Friederici and Wessels 1993; Jusczyk et al. 1994; Mattys and Jusczyk 2001). This may be achieved via learning the likelihood of phoneme transitions: for example, in English, certain phoneme transition

Received: June 22, 2023. Revised: March 21, 2024. Accepted: March 22, 2024

© The Author(s) 2024. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

probabilities are statistically unlikely (or even nonexistent, e.g. /dla/), while others are statistically more likely (e.g. /gla/). A similar principle is thought to be employed for syllable transitions, where statistically improbable syllable transitions can indicate between-word boundaries (Saffran et al. 1996).

Thus, while the initial analysis of phonemes is based on their acoustic features (Mesgarani et al. 2014; Khalighinejad et al. 2017; Overath and Lee 2017; Yi et al. 2019), subsequent processing stages are likely more linguistic in nature, such as those identifying language-specific phonemes or phonotactics, or even higher-level processes underlying the analysis of syntax, semantics, or lexical access (Kutas and Hillyard 1983; Friederici et al. 1993; Kocagoncu et al. 2017). In the current study, we refer to the initial process as the “acoustic analysis” of phonemes and to the subsequent process as the “linguistic analysis” of phonemes. While decades of research have resulted in detailed speech/language models (Hickok and Poeppel 2007; Rauschecker and Scott 2009; Friederici 2011), a clear demarcation between acoustic and linguistic analyses (and their potential bi-directional interaction) that contribute toward speech comprehension has largely remained elusive. One reason for this is that, in everyday listening situations, acoustic and linguistic analyses are difficult to separate and likely interact, e.g. via top-down modulation of acoustic feature analysis by linguistic processes (Anderson et al. 2003; Davis and Johnsrude 2007; Díaz et al. 2008). In addition, previous studies that investigated phoneme processing in naturalistic contexts (Mesgarani et al. 2014; Khalighinejad et al. 2017; Daube et al. 2019; Gwilliams et al. 2022) did so only in a familiar language: this approach is unable to dissociate the initial acoustic processes from the obligatory nature of linguistic processes that become engaged in a native, familiar language.

In contrast, Overath and Lee (2017) were recently able to dissociate the acoustic and linguistic processes underlying phoneme analysis by comparing PRPs in familiar vs. foreign languages. They used a variant of a novel sound quilting algorithm (Overath et al. 2015) to create “speech-based quilts” in which linguistic units (phoneme, syllable, word) were pseudo-randomly “stitched together,” or quilted to form a new stimulus. (We introduce the term “speech-based quilting” here to emphasize the fact that the duration of stitched segments depends on the duration of each individual linguistic unit [here, phonemes of different duration]; in contrast, in the previously employed “time-based quilting” approach (Overath et al. 2015; Overath and Paik 2021), the segment duration is fixed [in ms].) This paradigm allowed the comparison of an acoustic stimulus manipulation (speech-based quilting) in a familiar vs. foreign language: if the processing of phonemes is affected by the acoustic manipulation (increasing linguistic unit size of speech quilts) in a familiar language only, then this would suggest that linguistic analysis in the familiar language influenced the acoustic analysis of phonemes. Put differently, if only minimal or no phonemic repertoire or phonotactic rules are available to a listener (as is the case in a foreign language), the encoding of speech sounds themselves would be independent of their ordering (phonotactics) or linguistic unit size in which they appear. Using EEG to investigate the PRP for different phoneme classes (Khalighinejad et al. 2017), Overath and Lee (2017) found that vowels in particular are amenable to such top-down linguistic modulation. However, the limited spatial resolution of EEG did not allow inferences as to where in the auditory cortex (or beyond) such top-down modulation might originate, or act upon.

Recent advances in functional magnetic resonance imaging (fMRI) time-series analysis have demonstrated that the neural

activity to natural speech stimuli can be predicted from fast-paced acoustic (e.g. envelope, spectrum), phonological, and semantic features via linearized encoding modeling (Huth et al. 2016; De Heer et al. 2017). Inspired by this approach, the current study employed linearized encoding modeling of fMRI data in human cortex in an effort to reveal the separate encoding of acoustic and linguistic features of speech. Specifically, we used speech-based quilting (original speech vs. phoneme quilts) in familiar (English) vs. foreign (Korean) languages to dissociate the neural correlates of the acoustic and linguistic processes that contribute to the analysis of a fundamental linguistic unit, the phoneme. Importantly, as opposed to previous research from our group where the mean blood-oxygen-level-dependent (BOLD) levels of comparatively short stimuli were compared, the design of the current study allowed us to analyze prolonged time-locked dynamics in BOLD timeseries based on acoustic and linguistic features of naturalistic and manipulated speech stimuli. In addition, the novel quilting variant deliberately kept phonemic information largely intact, while previous studies (Overath et al. 2015; Overath and Paik 2021) largely destroyed such information by placing segment boundaries at strictly regular time intervals (regardless of phonemic boundaries). This allows for a more precise understanding of the temporal dynamics underlying the transformation of acoustic to linguistic information. We show (i) that the acoustic analysis of phonemes is modulated by linguistic processes and (ii) that the interaction cannot be explained by solely acoustic or phonetic information.

Materials and methods

Overview

Ten native English speakers without any knowledge of Korean listened to speech stimuli in 4 conditions (original speech or phoneme quilts, in English or Korean) during 3 sessions of fMRI scanning. Condition-specific linearized encoding models were trained to predict the fMRI time-series using 2 “Acoustic” predictors (the broadband envelope and its first-order derivative with positive half-wave rectification) and 4 “Phonetic” predictors (the durations of each of the four main phoneme classes; i.e. vowels, nasals and approximants, plosives, fricatives and affricatives). Note that Multipenalty ridge regression models were optimized using a Bayesian optimizer for each predictor (i.e. 6 regularization hyperparameters per model). The prediction accuracy of full models with all predictors was calculated using Pearson’s correlation (r). The unique contribution of a certain predictor group (or a feature subspace, e.g. Daube et al. (2019)), was calculated using partial correlation (ρ) by regressing out the other predictor group (see Fig. 1 for an overview of the analysis).

The raw data supporting the conclusions of this article will be made available by the authors upon request. In-house functions and scripts (MATLAB) to reproduce results in the manuscript as well as exemplar stimuli are publicly available on the Open Science Framework repository (<https://osf.io/zgj3m/>).

Participants

Ten native English speakers without any knowledge or experience in Korean participated in the current study (mean age = 24.0 ± 2.2 yr; 6 females). Eight participants volunteered in three sessions consisting of 8 runs each on separate days (intervals in days: mean = 8.5, standard deviation = 16.6, min = 1, max = 70) and two other participants in a single session each (6 and 8 runs, respectively), resulting in a total of 24 scanning sessions. This is on par with similar approaches that maximize

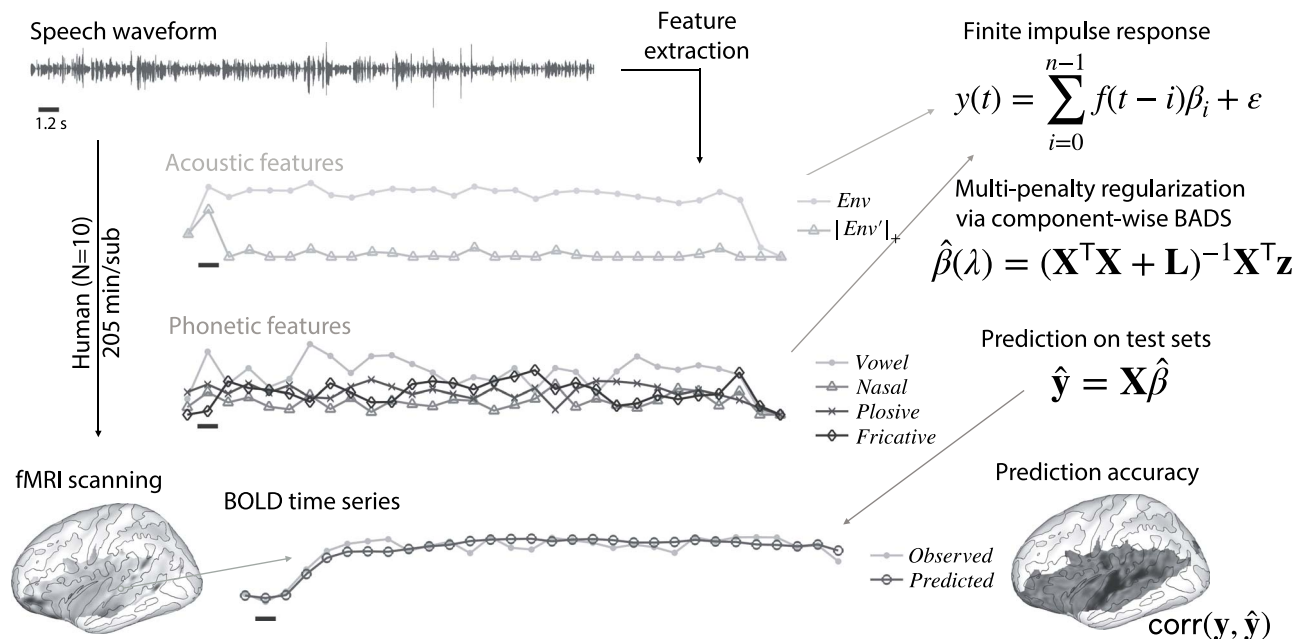


Fig. 1. Linearized encoding analysis overview. Functional MRI data were acquired from 10 human participants while listening to unmanipulated or phoneme-scrambled speech stimuli in either English or Korean. From the speech waveform, acoustic features (cochleogram envelope and its first-order derivative with positive half-wave rectification) and phonetic features (the duration of four phoneme classes) were extracted and down-sampled at the fMRI sampling rate (1/1.2 Hz). Scale bars represent an fMRI sampling period (1.2 s). After preprocessing, the surface-mapped BOLD time series $y(t)$ was predicted using regularized finite impulse response modeling. Multi-penalty regularization was optimized in a principal component space using the BADS optimizer. The cross-validated prediction accuracy was measured by Pearson correlation between observed and predicted BOLD time series that was back-projected onto the vertex space (see linearized encoding analysis section for details).

intrasubject reliability over intrasubject variability in the data (Kay et al. 2008; Moerel et al. 2013; Naselaris et al. 2015; Norman-Haignere et al. 2015; Huth et al. 2016; Santoro et al. 2017; Breedlove et al. 2020).

All participants were recruited via the Brain Imaging and Analysis Center at Duke University Medical Center, NC, USA after safety screening for MRI (e.g. free of metal implants and claustrophobia). All reported to have normal hearing and no history or presence of neurological or psychiatric disorders. Informed written consent was obtained from all participants prior to the study in compliance with the protocols approved by the Duke University Health System Institutional Review Board.

Stimuli

Speech stimuli were created from recordings (44,100 Hz sampling rate, 16-bit precision) of four female bilingual (Korean and English) speakers reading textbooks in either language as in previous studies (Overath and Lee 2017; Overath and Paik 2021). Native English and Korean speakers judged the recordings as coming from native English and Korean speakers, respectively. Korean was chosen because of its dissimilarity to English: it shares no etymological roots with English and has different syntactic and phonetic structures (Sohn 2001).

We used a modified version of the quilting algorithm (Overath and Lee 2017) where we pseudorandomized the order of phonemes (instead of set segment lengths). First, phonemes were extracted from the recordings and corresponding transcripts using the Penn Phonetic Lab Forced Aligner (https://babel.ling.upenn.edu/phonetics/old_website_2015/p2fa/index.html) (Yuan and Liberman 2008) for English speech and the Korean Phonetic Aligner (<https://korean.utsc.utoronto.ca/kpa/>; Yoon and Kang 2013) for Korean speech. The phoneme segmentation output was a Praat TextGrid, which was then imported to MATLAB

(<https://github.com/bbTomas/mPraat>) via the mPraat toolbox (Bořil and Skarnitzl 2016). The alignment was manually validated by a native English and Korean speaker, respectively (Overath and Lee 2017; Overath and Paik 2021). The durations of phonemes in the recordings of natural speech in milliseconds were as follows (see Supplementary Fig. S5a for histograms): min=4.3, max=396.2, mean=72.8, median=63.8, standard deviation=41.7, skewness=1.2 in English ($n=10,514$); min=8.9, max=308.3, mean=71.9, median=63.7, standard deviation=36.0, skewness=1.3 in Korean ($n=10,894$). The average durations were similar between languages (0.9-ms longer in English, $t[21406]=1.67$, $P=0.094$), while the distributions were slightly different for that English had more instances of short (e.g. < 20 ms) phonemes (Kolmogorov–Smirnov statistic=0.1413, $P=10^{-93}$).

The phoneme segments were pseudorandomly rearranged to create novel phoneme quilts. For each stimulus, a random initial phoneme was chosen; subsequent phonemes were selected such that the acoustic change at the boundary was as close as possible to the acoustic change in the original source signal (using the L2-norm metric of an equivalent rectangular bandwidth [ERB]-spaced cochleogram; see Overath et al. 2015). In addition, we applied the following exclusion criteria: (i) the phoneme duration needed to be at least 20 ms, (ii) two identical phonemes could not occur next to each other, and (iii) for a given phoneme, its subsequent phoneme could not be the same as in the original source signal. We used the pitch-synchronous overlap-add algorithm (Moulines and Charpentier 1990) to further minimize abrupt changes in pitch at phoneme boundaries. Overall alterations due to the quilting algorithm were quantified by the Kullback–Leibler divergence (D_{KL}) between L2-norm acoustic change distributions in the original source and the created phoneme quilt (median $D_{KL}=0.6873$ bits for English, 0.6004 bits

for Korean; Wilcoxon rank sum equal median test: $Z=0.5913$, $P=0.5543$). In the phoneme quilts, the durations of phonemes in milliseconds were as follows (Supplementary Fig. S5b): min = 20.0, max = 351.0, mean = 72.3, median = 63.0, standard deviation = 39.4, skewness = 1.4 in English ($n=10,467$); min = 20.0, max = 383.0, mean = 69.7, median = 60.0, standard deviation = 36.3, and skewness = 1.5 in Korean ($n=11,213$). There were slight differences between languages in average durations (2.6-ms longer in English, $t[21678]=5.08$, $P=10^{-7}$) and distributions (KS-stat = 0.0657, $P=10^{-21}$); however, the mean difference of 2.6 ms is much shorter than the modeled cochlear integration time-window of 20 ms.

As for the temporal modulation of the speech stimuli (Supplementary Fig. S3), the phoneme-based quilting decreased the temporal modulation energy at around 3–10 Hz in both languages. However, this reduction was greater in Korean than in English (max $F[1, 23]=41.09$; min FDR- $P=0.0005$; significant [FDR- $P < 0.05$] frequency bins = 3.9, 5.0, 5.1, 5.7, 6.4, 6.5, 6.6 Hz, in total 7 frequency bins; max effect size = 5.11 dB). This motivated us to include acoustic predictors in the encoding analysis (see Predictors section).

For both languages (English and Korean), the 33-s-long stimuli in the 2 experimental conditions (Original and Phoneme Quilts) were created by concatenating six 5.5-s stimuli (24 unique exemplars per condition and language) without gaps. Subsequent 5.5-s stimuli were either from the same or a different speaker (participants were asked to detect changes in the speaker, see Experimental procedure section). The overall sound intensity was normalized by equalizing the root-mean-square signal intensity across stimuli. At the beginning and the end of the 33-s stimuli, 10-ms cosine ramps were applied to avoid abrupt intensity changes.

Experimental procedure

Functional MRI data were acquired while participants listened to the speech stimuli (either Original or Phoneme Quilts in either language) and performed a task to maintain attention to the stimuli. A 33-s trial consisted of six 5.5-s stimuli of multiple speakers in a given condition. Silent intertrial intervals were uniformly varied between 5.6 and 10.4 s (mean = 8 s). One run consisted of twelve 33-s trials, and one session consisted of eight 8.5-min runs (except for 1 participant, who only completed 6 runs). For one of the 8 participants with 3 sessions, one run was prematurely terminated after 9 of the 12 trials due to technical difficulties (the intact 9 trials from the run were still used in the analysis). In total, fMRI data corresponding to ~203 min/participant were obtained for the 8 participants with 3 full sessions (average of ~174 min/participant for all 10 participants); this corresponds to ~158 min of stimuli (excluding the ISI) per participant with 3 full sessions (average of ~137 min/participant for all 10 participants).

The stimulus presentation timing was controlled via the Psychophysics Toolbox (v3.0.11 [<http://psycho toolbox.org/>]). Each run was triggered by the transistor-transistor-logic (TTL) signal from the MRI scanner mediated by a counter. Digital auditory signals at 44,100-Hz sampling rate and 16-bit precision from a Windows desktop were converted to analog signals by an external digital amplifier (Sony, Tokyo, Japan) and delivered to participants via MRI-compatible insert earphones (S14, Sensimetrics, MA, USA) at a comfortable listening level (~75 dB SPL). Participants wore protective earmuffs on top of the earphones to further reduce acoustic noise from the MRI scanner.

The task was to indicate a change in speaker (i.e. a 5.5-s stimulus of one speaker followed by a different speaker) via a button press on an MRI-compatible 4-button pad. Importantly, the participants could not know how many speaker changes occur as the number of changes also varied across 33-s trials (average speaker changes per trial = 3.5, between 1 and 4). The performance was assessed via d-prime $d' = \Phi^{-1}(\Pr(Y|s)) - \Phi^{-1}(\Pr(Y|n))$, where $\Pr(Y|s)$ is the hit rate in “signal” trials, $\Pr(Y|n)$ is the false alarm rate in “noise” trials, and $\Phi^{-1}(\bullet)$ is the inverse cumulative distribution function of the zero-mean, unit-variance Gaussian distribution (Macmillan and Kaplan 1985). Responses were classified as a hit if they occurred within 3 s following a change in speaker (and otherwise classified as false alarm). In the case of multiple responses within one 5.5-s stimulus segment, only the first response was counted. For extreme values of hit/false alarm rates (i.e. 0 or 1), an adjustment (i.e. adding $0.5/n$ to zero or subtracting $0.5/n$ from one for n trials) was made to avoid infinite values of d' (Macmillan and Kaplan 1985).

After each 33-s trial, participants received visual feedback about their performance ($D' = d' / \max d'$, where $\max d'$ is a d' for a perfect performance, ranging between [−100%, 100%]) with a description (“POOR” for $D' < 0$, “FAIR” for $0 \leq D' < 50\%$, “GOOD” for $50\% \leq D' < 100\%$, “PERFECT!” for $D' = 100\%$) to encourage continued attention. While multiple button presses were discarded from computing d' , an alerting message was presented to the participants (“NO KEY PRESSED!” or “TOO MANY KEYS PRESSED!”) instead of the performance feedback when the button presses were too many (>5) or none (2.5% of total 2397 trials from 9 participants; participant 1 was excluded from the d-prime analysis due to a technical fault of the in-scanner response device). The average D' was $61.1\% \pm 38.4\%$ points (overall: $d' = 1.14 \pm 0.72$; English-Original: 1.38 ± 0.23 ; English-Quilts: 1.02 ± 0.48 ; Korean-Original: 1.25 ± 0.38 ; Korean-Quilts: 0.95 ± 0.51). Repeated-measures ANOVA revealed a significant difference between original speech and phoneme quilts ($\eta_p^2 = 0.70$, $F[1,8]=16.37$, $P=0.02$; $d' = 1.28 \pm 0.62$ in original speech, $d' = 0.95 \pm 0.77$ in phoneme quilts) but neither between languages ($\eta_p^2 = 0.43$, $F[1, 8]=5.46$, $P=0.21$; $d' = 1.17 \pm 0.69$ in English, $d' = 1.05 \pm 0.74$ in Korean) nor an interaction between quilting and language ($\eta_p^2 = 0.09$, $F[1, 8]=0.68$, $P=0.43$).

Image acquisition

All images were acquired using a GE MR 750 3.0 Tesla scanner (General Electric, Milwaukee, WI, USA) with a 32-channel head coil system at the Duke University Hospital, NC, USA. For BOLD contrast, gradient-echo echo-planar imaging (GE-EPI) with a simultaneous multi-slice acceleration factor of 3 (i.e. 3 slices acquired in parallel with aliasing of FOV/3 shifts along the frequency-encoding direction) was used (in-plane pixel size = 2×2 mm², slice thickness = 2 mm, slice gap = 0 mm, FOV = 256×256 mm², matrix size = 128×128 , TE = 30 ms, flip angle = 73°, TR = 1200 ms, and phase-encoding direction = posterior-to-anterior). A total of 39 slices were acquired for each volume (13 slices per band) in an interleaved ascending sequence. At the beginning of a run, the volume was centered on the supratemporal plane, covering from the inferior colliculus to the inferior frontal gyrus. To correct for magnetic inhomogeneity artifacts, an additional GE-EPI image of 3 volumes with a reversed phase encoding direction (anterior-to-posterior) was acquired after each run except for the first participant.

For T_1 -weighted contrast, a magnetization prepared rapid gradient echo (MP-RAGE) scan covering the whole-brain (in-plane pixel size = 1×1 mm², slice thickness = 1 mm, slice gap = 0 mm,

FOV = 256 mm, matrix size = 256 × 256, TE = 3.2 ms, flip angle = 8°, TR = 2264 ms, and number of slices = 156) was acquired at the end of each session.

Image processing

Anatomical images

T₁-weighted images were segmented using SPM (SPM12; v7487 [<https://www.fil.ion.ucl.ac.uk/spm/>]) to obtain tissue probability maps (`spm.spatial.preproc`), which were used for anatomical CompCor regressors (Behzadi et al. 2007). High-resolution cortical surfaces were fully automatically constructed using FreeSurfer (v6.0.0 [<http://freesurfer.net/>]) for surface-based analysis.

Functional images

The displacement due to inhomogeneity in the B₀ field (i.e. susceptibility artifacts) was corrected using topup in FSL (v5.0.11 [<https://fsl.fmrib.ox.ac.uk/>]) with the reversed phase-encoding images. The first 6 volumes (i.e. “dummy scans”) of each run were subsequently discarded from the analyses. Temporal and spatial realignments were achieved using SPM: the slices were first temporally aligned to the center of the TR using sinc-interpolation (`spm.temporal.st`), and then, the volumes were spatially aligned to the mean volume using 4-th degree B-spline interpolation (`spm.spatial.realignunwarp`). Since we used a multiband sequence (i.e. 3 slices were acquired simultaneously), the acquisition time of each slice and reference time were provided (instead of slice order) for the slice-timing correction.

In order to suppress nonneural signal fluctuation, which is highly likely due to motion artifacts, we used the anatomical CompCor denoising technique (Behzadi et al. 2007). The specific steps were as follows: based on co-registered tissue segmentation probability maps from SPM, voxels with >99% tissue probability were selected. Subsequently, on concatenated time series from the voxels, principal component analysis (PCA) was applied to extract principal components. Among the extracted components, the top 6 components with the highest eigenvalues were used as motion regressors in the general linear model (GLM) denoising procedure (see Surface-based GLM denoising section).

Next, using FreeSurfer (Fischl 2012), the EPI volumes were projected onto individual cortical surfaces (~150,000 vertices per hemisphere) at the middle depth of cortices by averaging samples at the 40%, 50%, and 60% of cortical thickness to avoid aliasing (`mri_vol2surf` in FreeSurfer). Surface-mapped functional data were normalized to “`fsaverage6`” surfaces (40,962 vertices per hemisphere) via spherical surface registration and then smoothed with a 2D Gaussian kernel with the full-width-at-half-maximum of 6 mm (i.e. 3 pixels in the EPI slices) via iterative nearest-neighbor averaging (`mri_surf2surf` in FreeSurfer).

Surface-based GLM denoising

We applied a model-based denoising technique for task-based fMRI data (GLMdenoise v1.4 [<https://kendrickkay.net/GLMdenoise/>]) to the surface-mapped data (Kay et al. 2013). The algorithm extracts “noise” regressors from the data that would increase prediction accuracy in leave-one-run-out-cross-validation. This is achieved by first defining “noise pool” vertices with negative R² values for a given design matrix (i.e. vertices that are irrelevant to the task of interest), extracting principal components from the noise pool, and then determining an optimal number of components to remove as a minimal number where the improvement in cross-validation (CV) prediction

decays. We used box-car functions to represent the 4 conditions in the design matrix. On average, 4.5 ± 2.1 noise regressors were regressed out. These improved reliability in estimation (mean over standard errors ratio of coefficients estimates across CV folds: median increase = 0.82; mean increase = 1.12) but only slightly increased prediction accuracy (CV R²: median increase = 0.25% points; mean increase = 0.56% points). In addition to the noise regressors, the 4th order polynomial fits to slow drifts in BOLD time series, the 6 CompCor regressors, and the button-press regressors convoluted with a canonical HRF were regressed out from the residuals (i.e. prediction from the design matrix subtracted from the data).

Linearized encoding analysis

We predicted BOLD time series at each vertex in response to speech sounds using a linearized encoding model based on finite-impulse response (FIR) functions. Multiple lags were used to model the variable hemodynamic responses in different cortical areas (Huth et al. 2016; De Heer et al. 2017). In order to account for the collinearity of predictors representing acoustic and phonetic information, we used ridge regression to fit the model (i.e. FIR weights) and evaluated the prediction via CV. The procedures are explained in detail in the following subsections.

Vertex selection

For our interest in auditory and linguistic processing, we restricted our analysis to vertices in cortical regions that are previously known to be involved in speech processing so as to avoid unnecessary computations. Specifically, from the automatic parcellation based on the Desikan–Killiany cortical atlas (Desikan et al. 2006), the following 19 labels were included: “bankssts,” “caudalmiddlefrontal,” “inferiorparietal,” “inferiortemporal,” “lateralorbitofrontal,” “middletemporal,” “parsopercularis,” “parsorbitalis,” “parstriangularis,” “postcentral,” “precentral,” “rostralmiddlefrontal,” “superiorparietal,” “superiortemporal,” “supramarginal,” “frontalpole,” “temporalpole,” “transversetemporal,” and “insula.” The regions of interest are visualized in [Supplementary Fig. S6](#). Out of the vertices in the regions of interest, vertices with BOLD time series (i.e. where the acquisition slices of the EPI sequence were positioned) were individually selected. These vertices slightly varied across participants due to the variability of head sizes, individual acquisition volumes at each session, and movements across runs during sessions. [Supplementary Figure S7](#) shows the overlap of selected vertices across participants. On average, 28,297 ± 3,748 vertices were selected per participant. Note that, for the group-level statistical analysis, only vertices with data from all participants were included (18,818 vertices across both hemispheres).

Predictors

We included as predictors (i) the durations of phoneme classes (vowels, nasals and approximants, plosives, fricatives and affricatives; Vo, Na, Pl, Fr, respectively) and (ii) the speech envelope and its first-order derivative with positive half-wave rectification. For (i), the onset time and duration of each phoneme were determined and then grouped according to phoneme class (Ladefoged and Johnstone 2015; Shin 2015) (see [Supplementary Table S1](#)). Bigram transition probabilities between phoneme classes ([Supplementary Fig. S8](#)) were effectively altered by the quilting algorithm (Hotelling’s T^2 between Original and Phoneme-quilts = 1563, $P < 10^{-6}$ for English; Hotelling’s T^2 = 1258, $P < 10^{-6}$ for Korean). The durations of phoneme classes were

modeled as box-car functions at the audio sampling rate (44.1 kHz) and were then down-sampled to $1/\text{TR}$ ($1/1.2=0.833$ Hz) following anti-aliasing low-pass filtering. To align with the slice timing correction applied to the BOLD time series, the resampled time points were also at the center of the TR. For (ii), the speech envelope was computed from a cochleogram (30 filters from 20 to 10,000 Hz, equally spaced on an ERB scale) by raising the Hilbert envelope of the resulting cochleogram to a power of 0.3 to simulate cochlear compression and summing energy across all 30 ERB channels (McDermott and Simoncelli 2011; Overath et al. 2015). The speech envelope was then down-sampled as for the phoneme class durations. The rectified derivative was calculated following the down-sampling to reflect slow temporal modulations. It is important to clarify that the speech envelope and its derivative were used as parsimonious models to represent acoustic energy and temporal modulation (Daube et al. 2019), without suggesting that they represent “purely acoustic” characteristics. Due to the inherent covariance structure in natural speech, there is inevitably some overlap between acoustic and linguistic information. This overlap is particularly pronounced in familiar languages. However, for foreign languages, such as Korean speech to native English speakers as in the current study, a separation of acoustic and linguistic information is feasible. This has been demonstrated in a series of studies conducted by our group (Overath and Lee 2017; Overath and Paik 2021).

The down-sampled predictors showed moderate collinearity; the diagnostic metric of the collinearity, namely “condition index,” had an overall value of 27 (23–24 in English conditions and 41 in Korean conditions) given a suggested criterion (>30) for a moderate multicollinearity (Belsley 1991). The condition index is a square root of the maximum eigenvalue divided by the minimum eigenvalue of the design matrix, quantifying the upper bound of the collinearity of the design matrix. The observed multicollinearity was mainly due to the high dependency between the vowel and the envelope predictors; the proportions of explained variance by the corresponding eigenvector (i.e. variance decomposition proportion) were 0.87 and 0.99 for the vowels and the envelope, respectively. The collinearity patterns were similar across conditions (Supplementary Fig. S9). The existence of multicollinearity motivated the use of a penalized regression.

Regularized FIR modeling

A FIR model was used to predict the BOLD time series at each vertex. In this approach, we modeled the neural response as a convolution of the predictors and a linear FIR filter, which is a commonly used approach in receptive field mapping of neural populations (Ringach et al. 1997; Wu et al. 2006).

Consider a linear model for t time points and p predictors,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1)$$

where \mathbf{y} is a $(t \times 1)$ data vector (i.e. BOLD time series at a certain vertex), \mathbf{X} is a $(t \times p)$ design matrix (i.e. a FIR model), $\boldsymbol{\beta}$ is a $(p \times 1)$ unknown coefficient vector, and $\boldsymbol{\varepsilon}$ is a noise vector from a zero-mean Gaussian distribution with a serial correlation $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \boldsymbol{\Omega})$, where $\boldsymbol{\Omega}$ is a $(t \times t)$ unknown covariance matrix and σ^2 is a scale factor. For the FIR modeling, the design matrix \mathbf{X} consists of matrices of delayed features as

$$\mathbf{X} = \begin{bmatrix} f_1 & f_2 & \cdots & f_p \end{bmatrix} * \mathbf{H}(n),$$

for p features and n delays as implemented in a convolutional kernel $\mathbf{H}(n)$, while $*$ denotes the convolution operation. A Toeplitz matrix can be constructed for delayed features between time point t_1 and t_2 with n delays for the i th feature as

$$f_i(t_1, t_2) * \mathbf{H}(n) = \begin{bmatrix} f_i(t_1) & f_i(t_1-1) & \cdots & f_i(t_1-(n-1)) \\ f_i(t_1+1) & f_i(t_1) & \cdots & f_i(t_1-n) \\ \vdots & \vdots & \cdots & \vdots \\ f_i(t_2) & f_i(t_2-1) & \cdots & f_i(t_2-(n-1)) \end{bmatrix},$$

where $f_i(t)$ is the scalar value of the i th predictor at time point t . In the current study, we delayed the predictors by 1, \dots , 10 TRs (1.2, \dots , 12 s). Once unknown coefficients (or weights) are estimated, an inner product $\mathbf{X}\hat{\boldsymbol{\beta}}$ is effectively a convolution of the i th feature and the estimated filter.

While it is standard to pre-whiten the data when modeling autocorrelated noise for a Generalized Least Squares (GLS) solution (Aitken 1936), here we did not pre-whiten the model. This is because even with autocorrelated noise, an Ordinary Least Squares (OLSs) solution is still an unbiased estimator (only its efficiency is suboptimal) and because our goal was to estimate (predict) responses, not to infer significance. In particular, for the current data, GLS often yielded worse CV prediction than OLS. Therefore, we empirically determined not to pre-whiten the model.

As we detected a strong collinearity among the predictors, we applied L_2 -norm regularization to solve Equation (1), which is known as a multipenalty ridge solution (Hoerl and Kennard 1970)

$$\hat{\boldsymbol{\beta}}(\lambda) = (\mathbf{X}^T \mathbf{X} + \mathbf{L})^{-1} \mathbf{X}^T \mathbf{y}, \quad (2)$$

where $\hat{\boldsymbol{\beta}}(\lambda)$ is a vector of penalized estimates and \mathbf{L} is a regularization matrix as

$$\mathbf{L} = \begin{bmatrix} \lambda_1 \mathbf{I} & \mathbf{O} & & \mathbf{O} \\ \mathbf{O} & \lambda_2 \mathbf{I} & & \mathbf{O} \\ & & \ddots & \\ \mathbf{O} & & & \lambda_p \mathbf{I} \end{bmatrix},$$

with λ_i is a scalar regularization hyperparameter for the i th feature, \mathbf{I}_i is the $(n \times n)$ identity matrix, and \mathbf{O} is a zero matrix in appropriate dimensions. The multipenalty ridge has been recently re-introduced to the neuroscience community as “banded ridge” (Nunez-Elizalde et al. 2019).

Component-wise optimization and vertex-wise evaluation

We used nested CV to optimize and evaluate the models (Hastie et al. 2009). To avoid information leakage driven by the stimulus-evoked responses (Hasson et al. 2010; Kaufman et al. 2012), trials were partitioned in a way that there is no overlap of stimuli across partitions. The whole data were partitioned into 2 outer-CV folds (training set and test set) to evaluate model performance on unseen data; the outer-CV training set is further partitioned into 2 inner-CV folds (training set and validation set) to optimize hyperparameters on independent data (Varoquaux et al. 2017).

Optimizing multiple penalty terms can be a nontrivial task (van de Wiel et al. 2021). Grid search algorithm can be efficiently implemented using GLM: i.e. a single inversion of the regularized design covariance matrix $(\mathbf{X}^T \mathbf{X} + \mathbf{L})$ can be used for all vertex-wise models for a particular search value. However, with the increasing number of penalty terms, the combinations of search

values exponentially increase. Bayesian adaptive direct search (BADs) optimizer [<https://github.com/lacerbi/bads>] is known to be robust in optimizing high-dimensional parameters (Acerbi and Ma 2017). But, as the BADs algorithm finds a unique optimization path in the parameter space, it runs through unique combinations of search values for each initialization. That is, separate inversions of differently regularized covariance matrices are required for each model; this process cannot be done simultaneously across all vertices (models) unlike the grid search using GLM.

To reduce the number of models to optimize, we exploited the spatial dependency of fMRI and data dimensionality (i.e. much fewer time-points than vertices). Using PCA on temporally concatenated data, we transformed the vertex time-series into the component space as

$$\mathbf{Z} = \mathbf{Y}\mathbf{A},$$

where \mathbf{Z} is a $(u \times k)$ component time-series matrix, \mathbf{Y} is a $(u \times v)$ vertex time-series matrix for u time-points (with all trials are temporally concatenated) and v vertices; $28,297 \pm 3,748$ (on average), and \mathbf{A} is a $(v \times k)$ demixing matrix for k components, which was determined for each subject to explain 99% of the total variance of the vertex time-series ($1,486 \pm 152$ on average; 5.25% of the number of vertices). Instead of a vertex-wise model (Equation 2), we fitted a component-wise model as

$$\hat{\gamma}(\lambda) = (\mathbf{X}^T \mathbf{X} + \mathbf{L})^{-1} \mathbf{X}^T \mathbf{z}, \quad (3)$$

where γ is a coefficient vector for a component and \mathbf{z} is a component time-series for a given component. The BADs optimizer searched the exponents of base 10 for λ s with absolute bounds of $[-15, 15]$ and plausible bounds of $[-10, 10]$. Over 10 random initializations (uniformly sampled within the plausible bounds), a vector of exponents that minimizes the validation error (sum of squares) was chosen for each component. The optimization process took 79–102 h per outer-CV fold, depending on the number of sessions, on Intel Xeon Gold 6130 processors (32 threads).

For evaluation, we predicted the component time-series for the outer-CV test set using the weights estimated from the outer-CV training set: $\hat{\mathbf{z}}_{te} = \mathbf{X}_{te} \hat{\gamma}_{tr}(\mathbf{L}^*)$, where subscripts “tr” and “te” indicate the outer-CV training set and test set, respectively, and \mathbf{L}^* denotes the optimal regularization matrix. Then, the predicted component time-series was transformed back into the vertex space: $\hat{\mathbf{Y}}_{te} = \hat{\mathbf{z}}_{te} \mathbf{W}$, where $\mathbf{W} = \mathbf{A}^{-1}$ is a $(k \times v)$ mixing matrix. Finally, for each vertex, the prediction accuracy was calculated by Pearson’s correlation (r) between the observed time-series and the predicted time-series: $r = \text{corr}(\mathbf{y}_{te}, \hat{\mathbf{y}}_{te})$.

To determine the uniquely explained variance by a particular set (subspace) of predictors (either the Acoustic or Phonetic subspace), vertex time-series were separately predicted based on each subspace as

$$\begin{cases} \hat{\mathbf{Y}}_A = \mathbf{X}_A \hat{\gamma}_A(\mathbf{L}) \mathbf{W} \\ \hat{\mathbf{Y}}_P = \mathbf{X}_P \hat{\gamma}_P(\mathbf{L}) \mathbf{W} \end{cases},$$

where $\mathbf{X} = [\mathbf{X}_A \quad \mathbf{X}_P]$ and $\gamma = [\gamma_A \quad \gamma_P]$. Note that the sum of separate predictions is equal to the full model prediction $\hat{\mathbf{Y}} = \hat{\mathbf{Y}}_A + \hat{\mathbf{Y}}_P$. Then, partial correlation (ρ) was calculated as Pearson’s correlation between residuals after regressing out one prediction from the data and the other prediction as $\rho_A = \text{corr}(\mathbf{y}, \hat{\mathbf{y}}_A; \hat{\mathbf{y}}_P)$ and $\rho_P = \text{corr}(\mathbf{y}, \hat{\mathbf{y}}_P; \hat{\mathbf{y}}_A)$.

Statistical inference

Statistical inference was computed via a nonparametric paired t-test using a cluster-based permutation test at group-level (Maris

and Oostenveld 2007). Specifically, r values of both models were calculated for each participant ($n=10$), and then, the difference between two models at each vertex was calculated. Next, the signs of differences across participants were flipped over all possible permutations ($2^{10}=1,024$) to form a null distribution. One-tailed P -values were computed from the null distribution for our directed hypotheses (e.g. English > Korean, Original > Quilts). Note that the inference was computed at the group-level, not the subject-level. Vertex-wise multiple comparisons correction was applied via a cluster-based permutation test as implemented in `ft_statistics_montecarlo` in FieldTrip (v20180903) (<http://www.fieldtriptoolbox.org/>) with a custom modification of clusterstat for a faster cluster identification through parallelization. In an earlier fMRI methodological study (Eklund et al. 2016), it was shown that a liberal cluster-forming threshold (CFT) in a cluster-level inference based on the random field theory resulted in a severely inflated family-wise error rate (FWER), whereas the permutation test showed a consistently proper control of the FWER regardless of the choice of a CFT. A recent study formally showed that a CFT in permutation tests does not affect the FWER, but only the sensitivity (Maris 2019). Thus, in the current study, clusters were defined by an arbitrary threshold of the α -level of 0.05 (for vertex-wise P -values) to improve the sensitivity, and the cluster-wise P -values are thresholded at the α -level of 0.005 to control the FWER at 0.005.

Results

Linguistic processing interacts with acoustic processing in the left superior temporal sulcus

We first investigated whether the encoding models could replicate our previous findings of mean BOLD activity (Overath et al. 2015; Overath and Paik 2021). Figure 2 displays the differences in Pearson correlation of full models between conditions (see Supplementary Fig. S1 for a rendering on uninflated surfaces). The original speech stimuli evoked BOLD time-series that are better explained by the encoding models than the quilted stimuli in the superior temporal sulci (STSs) and the anterior superior temporal gyri, bilaterally (Fig. 2a; max $t[9]=11.22$, min cluster- $P < 0.001$, max cluster size = 1,078 vertices, max $\Delta r = 0.2302$). The native language (English) as compared with the foreign language (Korean) showed similar but larger clusters over the lateral convexity of the STG (i.e. Te3; Morosan et al. (2005)), extending to the planum temporale in the left hemisphere (Fig. 2b; max $t[9] = 10.12$, min cluster- $P < 0.001$, max cluster size = 1,697 vertices, max $\Delta r = 0.1913$). An interaction in the expected direction (i.e. a greater difference for [Original > Quilts] in English than in Korean) was found in the superior portion of the left STS (i.e. Te4; Fig. 2c; max $t[9] = 7.98$, min cluster- $P = 0.002$, max cluster size = 484 vertices, max $\Delta r = 0.1147$), suggesting that the change of neural encoding as a function of the acoustic context was modulated by linguistic knowledge in this area.

Neither acoustic nor phonetic predictors can exclusively explain the interaction between acoustic and linguistic processes

After establishing the interaction of acoustic and linguistic processes, we further examined whether the difference between original and quilted speech in the native language was driven by processes at the acoustic or phonetic level. An additional motivation for this analysis to separate acoustic and phonetic contributions was to account for the fact that the phoneme-based quilting procedure had slightly different effects on the envelope modulation spectrum in the two languages (Supplementary Fig. S3).

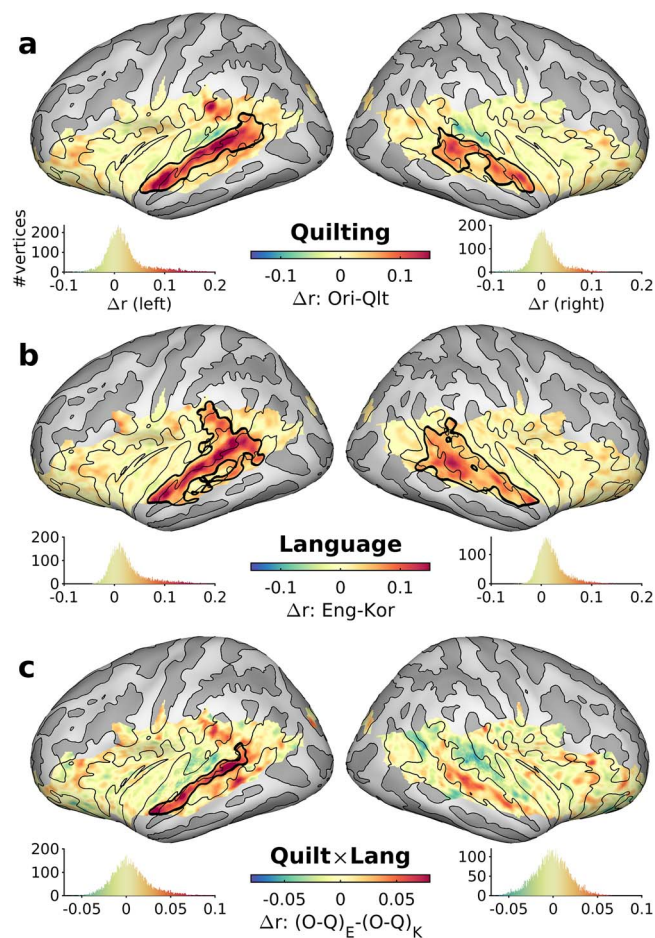


Fig. 2. Comparisons of prediction accuracies between conditions: a) main effect of quilting, original > quilts; b) main effect of language, English > Korean; and c) interaction between quilting and language, [English-original > English-quilts] > [Korean-original > Korean-quilts]. Thick black contours mark significant clusters (cluster- $P \leq 0.005$). Curvatures of the cortical surface are displayed in brighter (convex) and darker (concave) grays with thin black isocontours at the curvature of zero. Colored histograms of the r differences are displayed below each hemisphere. See Supplementary Fig. S2 for comparisons at the subject-level.

To this end, we calculated the partial correlation between the predicted and observed fMRI timeseries based on only one group of predictors (Acoustic or Phonetic) while regressing out the predicted timeseries based on the other group of predictors. The partial correlation thus quantifies the prediction based on the unique information in one predictor group in relation to the information that is common in both predictor groups. Finding an effect (of e.g. quilting) in the partial correlations would indicate that the linguistic modulation (i.e. the difference between original and quilted sounds in the native language as compared with the foreign language) is supported only by the unique characteristics of either Acoustic or Phonetic predictor groups that is orthogonal to common information of both predictor groups.

We tested whether the main effects and interaction identified in Fig. 2 can be explained by unique contributions of either of the predictor groups. For this, we averaged the partial correlation values in vertices within each cluster identified in Fig. 2. This ROI-based comparison reveals that only for the main effect of Language, the partial correlation of the Acoustic predictors was

significantly positive (max $t[9] = 5.60$, min $P < 0.001$, max diff $\rho = 0.0921$; Fig. 3b). That is, the residual STG/STS activity was better explained by the Acoustic predictors, but not by Phonetic predictors, for the native language (English) compared with the foreign language (Korean). For the main effect of Quilting (Fig. 3a) and the interaction (Fig. 3c), neither of the predictor groups showed significant positive partial correlation differences (min $P = 0.142$), suggesting that the interaction shown in the full models in Fig. 2 cannot be explained by unique information of either of the two predictor groups but is instead due to common information to both of them (see whole-cortex analyses in Supplementary Fig. S4).

Discussion

The phoneme is the fundamental linguistic unit that determines the meaning of words. We show that the four main phoneme classes, as well as broadband envelopes, are encoded in fMRI data acquired while listening to continuous speech signals. The acoustic processes underlying this phoneme analysis are modulated by linguistic analysis, whereby the acoustic manipulation (original speech vs. phoneme quilts) affected speech encoding more in a familiar language than in a foreign language. The results also revealed that this modulation cannot be explained uniquely by either acoustic or phonetic predictors.

Linguistic modulation of acoustic analysis of phonemes

Our primary aim was to dissociate acoustic from linguistic processes, which would enable us to determine their interaction, i.e. whether linguistic processes modulate the acoustic analysis of phonemes. To this end, we found that the acoustic manipulation (phoneme quilts vs. natural speech) had a larger effect on phoneme processing in a familiar language (English) than in a foreign language (Korean). Since the acoustic manipulation was the same for both languages, this suggests that the greater difference between acoustic contexts was due to linguistic processes becoming engaged in a familiar language (however, see also Contributions of acoustic and phonetic features section for further partitioning by predictor groups). Linguistic processes such as phonotactic, syntactic, as well as semantic analyses might therefore modulate the acoustic processing of phonemes, e.g. via hierarchical predictive coding or minimizing prediction errors through top-down modulation (Rao and Ballard 1999; Friston and Kiebel 2009). To our knowledge, this is the first demonstration of such linguistic modulation of a fundamental linguistic unit using fMRI. However, these results align well with Overath and Lee (2017), who found similar evidence for top-down linguistic modulation of phonemic analysis using a different recording modality (EEG).

Perhaps, the best-known example of the modulatory influence of linguistic information is that of phonemic restoration (Warren 1970; Samuel 1981). In phonemic restoration, a phoneme is still subjectively “perceived” even if it is masked or replaced completely by noise. This is often interpreted as an advantageous adaptation to speech perception in noisy environments, where it is common for interrupting or masking sounds to last only for a few tens or hundreds of milliseconds (i.e. on a temporal scale that is commensurate with that of phonemes). The top-down predictive nature of this phenomenon is further highlighted by the fact that, if the acoustic information is ambiguous, a “best guess” phoneme is perceived (Samuel 1987; Leonard et al. 2016). In fact, there is a wealth of evidence for such restorative processes in

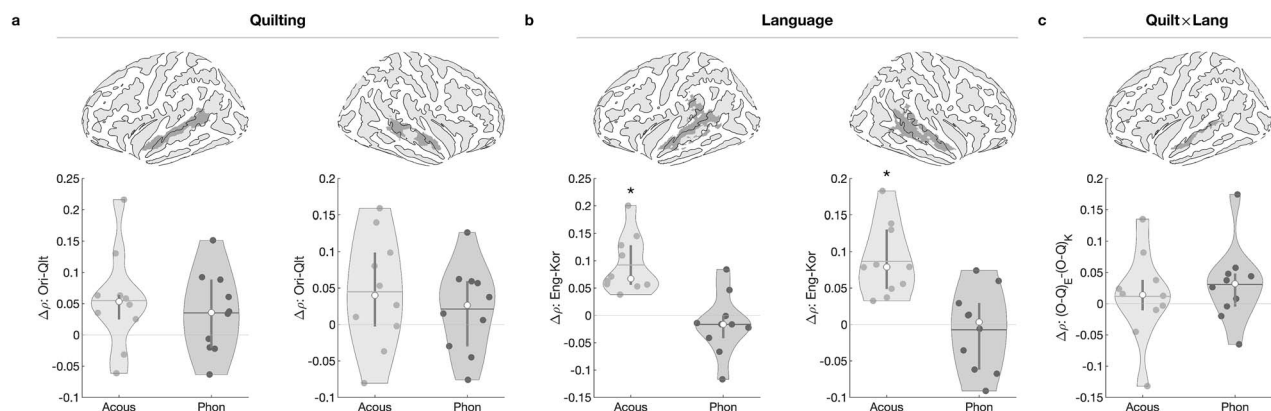


Fig. 3. ROI-based partial correlation differences for contrasts a) main effect of quilting (original > quilts), b) main effect of language (English > Korean), and c) interaction between quilting and language ([English-original > English-quilts] > [Korean-original > Korean-quilts]). On the cortical surface models, clusters are highlighted in dark gray. ROI-averaged partial correlations are shown in violin plots for the acoustic (light gray) and phonetic (dark gray) predictor groups; open circles mark medians, horizontal lines mark means, and vertical gray thick lines mark the first and third quartiles. *: Bonferroni-corrected $P < 0.005$.

speech perception, for example from studies using noise-vocoded stimuli (Shannon et al. 1995; Scott et al. 2000; Narain et al. 2003; Giraud et al. 2004; Obleser et al. 2008; Wild et al. 2012) or other methods to distort the speech signal (Davis et al. 2011; Eckert et al. 2016), while the most common explanation for restorative effects refers to top-down, predictive (Friston and Kiebel 2009) linguistic processes.

The locus of phonemic restoration, i.e. the region in which linguistic modulation is strongest, was recently shown to be situated in bilateral STG, likely due to receiving modulatory signals from left IFG (Leonard et al. 2016). This aligns remarkably well with the current study, where we found the strongest effect of linguistic modulation also along STG, albeit with a left-hemispheric dominance. The STG is a reasonable locus for such linguistic modulation, since it represents an intermediary processing stage in the language network that receives bottom-up information from primary auditory cortex and PT, as well as top-down information from higher-order auditory and frontal regions (Hickok and Poeppel 2007; Friederici 2009, 2011; Rauschecker and Scott 2009). For example, the analysis of spectral shape (a necessary computation to differentiate between the formant structures of different vowels) relies on bottom-up changes in effective connectivity between HG to PT, as well as PT to STG/STS regions (Warren et al. 2005; Kumar et al. 2007). In contrast, top-down signals from frontal cortex (e.g. left IFG) have been shown to modulate speech processing in auditory cortex (Sohoglu et al. 2012; Park et al. 2015; Cope et al. 2017; Overath and Paik 2021).

In the domain of electrophysiological measurements of speech perception, there is currently disagreement as to the extent that neural indices (such as speech-envelope entrainment, or phoneme encoding) can be interpreted as markers of linguistic processes that are necessary for speech comprehension (Luo and Poeppel 2007; Ding and Simon 2013; Di Liberto et al. 2015; Vanthornhout et al. 2018), or whether a more parsimonious explanation of these indices is that they reflect the analysis of characteristic acoustic properties of the speech signal (Howard and Poeppel 2010; Millman et al. 2015; Baltzell et al. 2017; Daube et al. 2019; Verschuere et al. 2021). Our study is able to shed new light on this controversy by directly comparing the encoding of acoustic properties of phonemes in either a familiar language or in a foreign language, in which no higher-level linguistic analysis takes place.

We should note that the current study did not measure linguistic processes explicitly. For example, participants did not perform a linguistic task (e.g. speech comprehension) but were simply asked to detect a change in speaker, a task that is largely orthogonal to linguistic processing (see also Overath and Paik (2021) for a similar task). Therefore, we interpret the linguistic modulation of phoneme class analysis as obligatory linguistic processes that become engaged as soon as familiar linguistic templates (e.g. phonotactics, syntax, lexicon, semantics) are detected in the signal. Future studies will need to determine whether, and to what extent, these obligatory linguistic processes for phoneme analysis are malleable to various tasks that engage specific linguistic processes. For example, the neural processing of acoustic features in speech sounds has been shown to be enhanced or sharpened if they are task-relevant, attended to vs. ignored, or primed (Mesgarani and Chang 2012; Holdgraf et al. 2016; Leonard et al. 2016; Rutten et al. 2019), and similar processes might become engaged for phoneme class encoding.

Contributions of acoustic and phonetic features

Another major advantage of the current study is the neural encoding analysis to delineate unique contributions of the overlapping information embedded in speech stimuli. The partial correlation analyses revealed a unique contribution of the acoustic features to the stronger neural encoding for the native language as compared with a foreign language (i.e. the main effect of language). Speech envelope (in particular onsets) encoding in the lateral STG has been shown from ECoG data with a native language (Ogania and Chang 2019). As shown in the whole-cortex analyses (Supplementary Fig. S4), this effect was strongest in the bilateral early auditory areas (HG/PT). This is in line with a recent finding that showed stronger speech envelope encoding in a native speaker group compared with a nonnative speaker group with a low proficiency (Liberto et al. 2021).

However, no unique contribution of the acoustic or phonetic features was found for either the main effect of quilting or the interaction of language and quilting. This suggests that the stronger encoding of all features for the original than quilted speech samples (Fig. 2b) and their interaction with language (i.e. more so for English than for Korean stimuli; Fig. 2c) are driven by the shared information between the acoustic and

phonetic predictors. A moderate multicollinearity between some predictors was indeed detected (see Statistical inference section), presumably due to the similarity between the speech envelope and vowel predictors. Note, though, that each predictor was individually optimized using Bayesian optimizer to avoid suboptimal regularization for individual predictors; that is, it is likely that all individual predictors were optimally regularized in the present analyses. Taken together, the current fMRI data suggest that the shared information between envelope and vowel predictors supported the modulation of the encoding strength in the acoustic and linguistic contexts.

Encoding of envelope and phoneme classes in the BOLD time series

One of our preliminary aims was to confirm that rapid acoustic and phonetic features can be shown to be encoded in a hemodynamic response that is approximately two orders of magnitude slower (tens of milliseconds vs. seconds). Encoding of these features had previously been demonstrated using electro-/magneto-physiological methods, which afford commensurate millisecond temporal resolution (Di Liberto et al. 2015; Khalighinejad et al. 2017; Yi et al. 2019; Brodbeck et al. 2022; Gwilliams et al. 2022; Heilbron et al. 2022). Nevertheless, the novel use of linearized ridge-regression modeling of fMRI BOLD signal time series was recently employed to successfully (and separably) reveal the encoding of acoustic and phonetic features in a familiar language: De Heer et al. (2017) collected fMRI data while presenting continuous, natural speech, and were able to reveal that the acoustic speech envelope predicted the BOLD time series best in HG, whereas articulatory phonetic features were predicted most accurately in higher-level auditory cortex such as STG. More recently, an fMRI study demonstrated a cortical hierarchy at a much longer time-scale (word-level) where the parameters (forecast distance and forecast depth) of a large language model (GPT-2) were mapped from the HG to STG, and then IFG (Caucheteux and King 2022). The current study is in broad agreement with these findings: while both the acoustic and phonetic features were encoded over the language network, the acoustic features showed greater encoding in the supratemporal regions (HG/PT).

More broadly, our study confirms that neural responses to rapid speech features, which are temporally integrated over several hundreds of milliseconds in the BOLD time series, can be revealed using linearized encoding modeling. Such models take advantage of the spatially separated functional organization of auditory cortex, for example with respect to prominent acoustic features such as frequency, spectro-temporal modulations, or spectral bandwidth (Rauschecker and Tian 2004; Saenz and Langers 2014; Santoro et al. 2014; Baumann et al. 2015; Moerel et al. 2018). This should encourage the future use of more naturalistic stimulus paradigms that allow the investigation of the complex dynamics of linguistic processes (Hamilton and Huth 2020), as well as other higher-order processes such as music perception (Kim 2022).

Modulation of acoustic and linguistic contexts

The analyses of the two factors Quilting and Language were motivated by previous studies that investigated the processing of temporal speech structure using segment-based speech quilting. In particular, these studies showed sensitivity in STS to temporal speech structure in either only a foreign language (Overath et al. 2015) or both native and foreign languages (Overath and Paik 2021), which is comparable to a main effect of Quilting here. In addition, activity in left IFG revealed an interaction between

Quilting and Language and increased as a function of temporal speech structure only in the familiar language (Overath and Paik 2021). In the current study, Quilting and Language both had greater prediction accuracies in left STS, while their interaction in the same area was due to larger prediction accuracy differences for the Original vs. Quilts contrast in English vs. Korean.

For successful speech comprehension, the temporal dynamics of speech necessitate analyses at multiple scales that are commensurate with the average durations of phonemes, syllables, words, sentences, etc. This temporal hierarchy is thought to be reflected in a cortical processing hierarchy in which the neuronal temporal window of integration (Theunissen and Miller 1995) increases from primary auditory cortex via nonprimary auditory cortex to frontal cortex (e.g. Lerner et al. 2011; though see Blank and Fedorenko 2020; Norman-Haignere et al. 2022 for a recent counterargument against such a hierarchy). The current results of greater prediction accuracy in STS as a function of Quilting largely support this view. A novel finding is the left-hemispheric lateralization. However, it is possible that this was driven by the interaction between Quilting and Language.

It is important to note that the segment-based quilting in previous studies disrupted the speech signal to a larger degree than the speech-based quilting employed here. The shortest segment length (30 ms) used in the previous studies, together with their placement irrespective of linguistic units, likely resulted in no phonemes being left fully intact in the resulting speech quilt. In contrast, the current speech-based quilting procedure preserved the phonemes (though likely still disrupted co-articulation cues).

Limitations and future directions

One potential limitation of this study is the possibility of participants attending more to intelligible stimuli (original speech in a familiar language) than unintelligible stimuli (all other conditions). While behavioral performance was better for original speech compared with quilted speech, it did not differ between languages or was affected differentially by language, which suggests that participants' task engagement was not differentially affected. In addition, it should be noted that this issue is not unique to the current study design, and is in fact a common challenge in speech studies that use unintelligible control stimuli. Here, we used an unfamiliar language to manipulate access to linguistic knowledge. While previous research has employed synthetic approaches to parametrically manipulate intelligibility by degrading spectral structures (Blessner 1972; Davis et al. 2005), slow temporal modulation (Ghitza 2012), or fine temporal structure (Lorenzi et al. 2006), these methods can introduce acoustic artifacts that distinguish them from natural human speech even at the acoustic level. In contrast, our study, along with previous work from our group, used speech recordings from bilingual speakers in 2 widely distant languages. This approach ensured that the natural acoustic features and speaker-specific acoustic characteristics remained highly similar between languages. We believe that unaltered foreign speech would have generated more attention than acoustically modified unintelligible nonspeech sounds. However, we acknowledge the possibility that the modulation of bottom-up attention from a familiar language may have influenced the current findings to some degree.

Another methodological limitation of the current study is the use of only two acoustic conditions (original and phoneme-quilts), which differs from previous studies conducted by our group. The rationale behind this design was to reduce sampling variance by minimizing the number of conditions while increasing the data points per condition. However, this design choice may

present a challenge when disentangling phonetic processing from higher-level linguistic processing beyond the phoneme-level (e.g. syllables, words, sentences), which may highly correlate with lower-level processing. While this was not the primary focus of the current study, it is important to emphasize the caveat of an encoding model, which entails that it does not necessarily rule out the possibility of unmodeled variables driving the observed prediction (Naselaris et al. 2011). However, this challenge is not applicable to the unfamiliar language, as listeners did not have access to higher-level linguistic features. That is, if the heightened prediction in the original condition was solely driven by unmodelled higher-level linguistic features, this effect would not have been observed in the unfamiliar language. On the contrary, the main effect of quilting was found both in English and Korean conditions when compared separately (Supplementary Fig. S10). This suggests that the current findings are unlikely to be driven solely by any unmodelled higher-level linguistic features. This is also in agreement with results reported in Overath et al. (2015) and Overath and Lee (2017), which showed an effect of increasing temporal structure (akin to the original speech vs. phoneme-quilt comparison here) in STS in unfamiliar languages.

The current study makes a number of predictions for future studies investigating the acousto-linguistic transformation of speech. We show evidence for linguistic modulation of a fundamental linguistic unit, the phoneme, in native English speakers when listening to English speech, but not when listening to a foreign language for which participants had no linguistic repertoire. Therefore, while it is unlikely that the current results are specific to English phonemes, future studies should confirm this interaction, for example in native Korean participants who have no knowledge of English. Similarly, people who are perfectly bilingual in English and Korean should show evidence for linguistic modulation in both languages as a function of quilting, while those for whom both languages are foreign should not. Alternatively, a modulation of this interaction by the proficiency level of English among Korean participants might be explored. However, it is important to note that examining between-subject effects like this would require a higher statistical power to ensure a reliable estimate (Marek et al. 2022).

In addition, the fact that the linguistic modulation of the acoustic speech signal operates at an intermediate stage of linguistic analysis likely reflects its significance: if linguistic modulation starts at the level of phonemes, its ability to impact a later word processing stage is conceivably greater than if linguistic modulation only started at the word processing stage. Given the highly predictive nature of speech processing (see Linguistic modulation of acoustic analysis of phonemes section above), such modulation might be particularly helpful in situations in which the speech signal is compromised (e.g. in noisy conditions such as in a restaurant or bar). People with hearing loss (e.g. presbycusis) are a clinical population that is known to struggle in such situations, even with the help of hearing aids (Moore 1996; Shinn-Cunningham and Best 2008). It is therefore possible that (at least) one reason for their exacerbated speech comprehension difficulties in noisy situations is that the linguistic modulation of phonemes has deteriorated, thereby reducing the effectiveness of predictive speech processes. A similar argument might be made for people suffering from “hidden hearing loss”: i.e. hearing difficulties without detectable deficits in routine audiometry tests (Kujawa and Liberman 2009; Ruggles et al. 2011). We predict that linguistic modulation of phoneme analysis is reduced in these populations (particularly in situations with background noise) and might thus serve as a clinical marker.

Conclusions

In conclusion, the current study demonstrates that individual phoneme classes derived from continuous speech signals are encoded in the BOLD signal time series. In particular, by using a design that dissociates acoustic from linguistic processes, we show that the acoustic processing of a fundamental linguistic unit, the phoneme, is modulated by linguistic analysis. The fact that this modulation operates at an intermediate stage likely enhances its ability to impact subsequent, higher-level processing stages, and as such might represent an important mechanism that facilitates speech comprehension in challenging listening situations.

Acknowledgments

All authors certify that they have no affiliations with or involvement in any organization or entity with any financial interest or non-financial interest in the subject matter or materials discussed in this article. The authors would like to thank Frankie Pennington and Joon Hyun Paik for manually checking phoneme onsets and offsets for the forced alignment of English and Korean recordings, respectively. The authors would also like to thank Dr. Mark A. van de Wiel for technical suggestion on multi-penalty ridge optimization.

Author contributions

SGK: Conceptualization, Investigation, Data curation, Methodology, Formal Analysis, Software, Visualization, Writing—original draft, Writing—review & editing; FDM: Formal Analysis, Methodology, Software, Writing—review & editing; TO: Funding acquisition, Project administration, Supervision, Resources, Conceptualization, Methodology, Writing—original draft, Writing—review & editing.

Supplementary material

Supplementary material is available at *Cerebral Cortex* online.

Funding

This work was supported by US National Institutes of Health grant R21DC016386 to TO. FDM has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. ERC - CoG 2020–101001270). The sponsor had no role in the study design, data collection, analysis, interpretation, preparation of the manuscript, or decision to submit the manuscript for publication.

Conflict of interest statement: None declared.

References

- Acerbi L, Ma WJ. Practical Bayesian optimization for model fitting with Bayesian adaptive direct search. In: Guyon I, Von Luxburg U, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R. (eds.). *Proceedings of Advances in Neural Information Processing Systems 30* Neural Information Processing Systems (La Jolla, CA), 2017:31:1837–1847.
- Aitken AC. On least squares and linear combination of observations. *Proc. R. Soc. Edinb. B.* 1936;55:42–48. <https://doi.org/10.1017/S0370164600014346>.

- Anderson JL, Morgan JL, White KS. A statistical basis for speech sound discrimination. *Lang Speech*. 2003;46(2-3):155–182. <https://doi.org/10.1177/00238309030460020601>.
- Baltzell LS, Srinivasan R, Richards VM. The effect of prior knowledge and intelligibility on the cortical entrainment response to speech. *J Neurophysiol*. 2017;118(6):3144–3151. <https://doi.org/10.1152/jn.00023.2017>.
- Baumann S, Joly O, Rees A, Petkov CI, Sun L, Thiele A, Griffiths TD. The topography of frequency and time representation in primate auditory cortices. *elife*. 2015;4:e03256. <https://doi.org/10.7554/eLife.03256>.
- Behzadi Y, Restom K, Liao J, Liu TT. A component based noise correction method (compcor) for bold and perfusion based fMRI. *NeuroImage*. 2007;37(1):90–101. <https://doi.org/10.1016/j.neuroimage.2007.04.042>.
- Belsley DA. A guide to using the collinearity diagnostics. *Médecine psychosomatique; regards sur les énigmes de la médecine*. 1991;4(1):33–50. <https://doi.org/10.1007/BF00426854>.
- Blank IA, Fedorenko E. No evidence for differences among language regions in their temporal receptive windows. *NeuroImage*. 2020;219:116925. <https://doi.org/10.1016/j.neuroimage.2020.116925>.
- Blessner B. Speech perception under conditions of spectral transformation. I. Phonetic characteristics. *J Speech Hear Res*. 1972;15(1):5–41. <https://doi.org/10.1044/jshr.1501.05>.
- Bořil T, Skarnitzl R. Tools rPraat and mPraat. In: Sojka P, Horák A, Kopeček I, Pala K (eds.). *Text, speech, and dialogue*. TSD 2016. *Lecture notes in computer science*. Cham: Springer International Publishing; 2016. pp. 367–374.
- Breedlove JL, St-Yves G, Olman CA, Naselaris T. Generative feedback explains distinct brain activity codes for seen and mental images. *Curr Biol*. 2020;30(12):2211–2224.e6. <https://doi.org/10.1016/j.cub.2020.04.014>.
- Brodbeck C, Bhattasali S, Cruz Heredia AAL, Resnik P, Simon JZ, Lau E. Parallel processing in speech perception with local and global representations of linguistic context. *elife*. 2022;11:e72056. <https://doi.org/10.7554/eLife.72056>.
- Caucheteux C, King J-R. Brains and algorithms partially converge in natural language processing. *Commun Biol*. 2022;5(1):134. <https://doi.org/10.1038/s42003-022-03036-1>.
- Cheour M, Ceponiene R, Lehtokoski A, Luuk A, Allik J, Alho K, Näätänen R. Development of language-specific phoneme representations in the infant brain. *Nat Neurosci*. 1998;1(5):351–353. <https://doi.org/10.1038/1561>.
- Chomsky N, Halle M. Some controversial questions in phonological theory. *J Linguist*. 1965;1(2):97–138. <https://doi.org/10.1017/S0022226700001134>.
- Cope TE, Sohoglu E, Sedley W, Patterson K, Jones PS, Wiggins J, Dawson C, Grube M, Carlyon RP, Griffiths TD, et al. Evidence for causal top-down frontal contributions to predictive processes in speech perception. *Nat Commun*. 2017;8(1):2154. <https://doi.org/10.1038/s41467-017-01958-7>.
- Daube C, Ince RAA, Gross J. Simple acoustic features can explain phoneme-based predictions of cortical responses to speech. *Curr Biol*. 2019;29(12):1924–1937.e9. <https://doi.org/10.1016/j.cub.2019.04.067>.
- Davis MH, Johnsrude IS. Hearing speech sounds: top-down influences on the interface between audition and speech perception. *Hear Res*. 2007;229(1-2):132–147. <https://doi.org/10.1016/j.heares.2007.01.014>.
- Davis MH, Johnsrude IS, Hervais-Adelman A, Taylor K, McGettigan C. Lexical information drives perceptual learning of distorted speech: evidence from the comprehension of noise-vocoded sentences. *J Exp Psychol Gen*. 2005;134(2):222–241. <https://doi.org/10.1037/0096-3445.134.2.222>.
- Davis MH, Ford MA, Kherif F, Johnsrude IS. Does semantic context benefit speech understanding through “top-down” processes? Evidence from time-resolved sparse fMRI. *J Cogn Neurosci*. 2011;23(12):3914–3932. https://doi.org/10.1162/jocn_a_00084.
- De Heer WA, Huth AG, Griffiths TL, Gallant JL, Theunissen FE. The hierarchical cortical organization of human speech processing. *J Neurosci*. 2017;37(27):6539–6557. <https://doi.org/10.1523/JNEUROSCI.3267-16.2017>.
- Desikan RS, Ségonne F, Fischl B, Quinn BT, Dickerson BC, Blacker D, Buckner RL, Dale AM, Maguire RP, Hyman BT, et al. An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage*. 2006;31(3):968–980. <https://doi.org/10.1016/j.neuroimage.2006.01.021>.
- Di Liberto GM, O’Sullivan JA, Lalor EC. Low-frequency cortical entrainment to speech reflects phoneme-level processing. *Curr Biol*. 2015;25(19):2457–2465. <https://doi.org/10.1016/j.cub.2015.08.030>.
- Díaz B, Baus C, Escera C, Costa A, Sebastián-Gallés N. Brain potentials to native phoneme discrimination reveal the origin of individual differences in learning the sounds of a second language. *Proc Natl Acad Sci USA*. 2008;105(42):16083–16088. <https://doi.org/10.1073/pnas.0805022105>.
- Ding N, Simon JZ. Adaptive temporal encoding leads to a background-insensitive cortical representation of speech. *J Neurosci*. 2013;33(13):5728–5735. <https://doi.org/10.1523/JNEUROSCI.5297-12.2013>.
- Eckert MA, Teubner-Rhodes S, Vaden KI Jr. Is listening in noise worth it? The neurobiology of speech recognition in challenging listening conditions. *Ear Hear*. 2016;37(1):101S–110S. <https://doi.org/10.1097/AUD.0000000000000300>.
- Eklund A, Nichols TE, Knutsson H. Cluster failure: why fMRI inferences for spatial extent have inflated false-positive rates. *Proc Natl Acad Sci USA*. 2016;113(28):7900–7905. <https://doi.org/10.1073/pnas.1602413113>.
- Fischl B. Freesurfer. *NeuroImage*. 2012;62(2):774–781. <https://doi.org/10.1016/j.neuroimage.2012.01.021>.
- Friederici AD. Pathways to language: Fiber tracts in the human brain. *Trends Cogn Sci*. 2009;13(4):175–181. <https://doi.org/10.1016/j.tics.2009.01.001>.
- Friederici AD. The brain basis of language processing: from structure to function. *Physiol Rev*. 2011;91(4):1357–1392. <https://doi.org/10.1152/physrev.00006.2011>.
- Friederici AD, Wessels JMI. Phonotactic knowledge of word boundaries and its use in infant speech perception. *Percept Psychophys*. 1993;54(3):287–295. <https://doi.org/10.3758/BF03205263>.
- Friederici AD, Pfeifer E, Hahne A. Event-related brain potentials during natural speech processing: effects of semantic, morphological and syntactic violations. *Cogn Brain Res*. 1993;1(3):183–192. [https://doi.org/10.1016/0926-6410\(93\)90026-2](https://doi.org/10.1016/0926-6410(93)90026-2).
- Friston K, Kiebel S. Predictive coding under the free-energy principle. *Philos Trans R Soc Lond Ser B Biol Sci*. 2009;364(1521):1211–1221. <https://doi.org/10.1098/rstb.2008.0300>.
- Ghitza O. On the role of theta-driven syllabic parsing in decoding speech: intelligibility of speech with a manipulated modulation Spectrum. *Front Psychol*. 2012;3:238. <https://doi.org/10.3389/fpsyg.2012.00238>.
- Giraud AL, Kell C, Thierfelder C, Sterzer P, Russ MO, Preibisch C, Kleinschmidt A. Contributions of sensory input, auditory search

- and verbal comprehension to cortical activity during speech processing. *Cereb Cortex*. 2004;14(3):247–255. <https://doi.org/10.1093/cercor/bhg124>.
- Gwilliams L, King J-R, Marantz A, Poeppel D. Neural dynamics of phoneme sequences reveal position-invariant code for content and order. *Nat Commun*. 2022;13(1):6606. <https://doi.org/10.1038/s41467-022-34326-1>.
- Hamilton LS, Huth AG. The revolution will not be controlled: natural stimuli in speech neuroscience. *Lang Cogn Neurosci*. 2020;35(5):573–582. <https://doi.org/10.1080/23273798.2018.1499946>.
- Hasson U, Skipper JI, Nusbaum HC, Small SL. Abstract coding of audiovisual speech: beyond sensory representation. *Neuron*. 2007;56(6):1116–1126. <https://doi.org/10.1016/j.neuron.2007.09.037>.
- Hasson U, Malach R, Heeger DJ. Reliability of cortical activity during natural stimulation. *Trends Cogn Sci*. 2010;14(1):40–48. <https://doi.org/10.1016/j.tics.2009.10.011>.
- Hastie T, Tibshirani R, Friedman J. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, New York, NY; 2009.
- Heilbron M, Armeni K, Schoffelen J-M, Hagoort P, de Lange FP. A hierarchy of linguistic predictions during natural language comprehension. *Proc Natl Acad Sci*. 2022;119(32):e2201968119. <https://doi.org/10.1073/pnas.2201968119>.
- Hickok G, Poeppel D. The cortical organization of speech processing. *Nat Rev Neurosci*. 2007;8(5):393–402. <https://doi.org/10.1038/nrn2113>.
- Hoerl AE, Kennard RW. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*. 1970;12(1):55–67. <https://doi.org/10.1080/00401706.1970.10488634>.
- Holdgraf CR, de Heer W, Pasley B, Rieger J, Crone N, Lin JJ, Knight RT, Theunissen FE. Rapid tuning shifts in human auditory cortex enhance speech intelligibility. *Nat Commun*. 2016;7(1):13654. <https://doi.org/10.1038/ncomms13654>.
- Howard MF, Poeppel D. Discrimination of speech stimuli based on neuronal response phase patterns depends on acoustics but not comprehension. *J Neurophysiol*. 2010;104(5):2500–2511. <https://doi.org/10.1152/jn.00251.2010>.
- Huth AG, de Heer WA, Griffiths TL, Theunissen FE, Gallant JL. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*. 2016;532(7600):453–458. <https://doi.org/10.1038/nature17637>.
- Jusczyk PW, Luce PA, Charles-Luce J. Infants' sensitivity to phonotactic patterns in the native language. *J Mem Lang*. 1994;33(5):630–645. <https://doi.org/10.1006/jmla.1994.1030>.
- Kaufman S, Rosset S, Perlich C, Stitelman O. Leakage in data mining: formulation, detection, and avoidance. *ACM Trans Knowl Discov Data*. 2012;6(4):1–21. <https://doi.org/10.1145/2382577.2382579>.
- Kay KN, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008;452(7185):352–355. <https://doi.org/10.1038/nature06713>.
- Kay K, Rokem A, Winawer J, Dougherty R, Wandell B. Glmde-noise: a fast, automated technique for denoising task-based fMRI data. *Front Neurosci*. 2013;7:247. <https://doi.org/10.3389/fnins.2013.00247>.
- Khalighinejad B, Cruzatto da Silva G, Mesgarani N. Dynamic encoding of acoustic features in neural responses to continuous speech. *J Neurosci*. 2017;37(8):2176–2185. <https://doi.org/10.1523/JNEUROSCI.2383-16.2017>.
- Kilian-Hütten N, Valente G, Vroomen J, Formisano E. Auditory cortex encodes the perceptual interpretation of ambiguous sound. *J Neurosci*. 2011;31(5):1715–1720. <https://doi.org/10.1523/JNEUROSCI.4572-10.2011>.
- Kim S-G. On the encoding of natural music in computational models and human brains. *Front Neurosci*. 2022;16:928841. <https://doi.org/10.3389/fnins.2022.928841>.
- Kleinschmidt DF, Jaeger TF. Robust speech perception: recognize the familiar, generalize to the similar, and adapt to the novel. *Psychol Rev*. 2015;122(2):148–203. <https://doi.org/10.1037/a0038695>.
- Kocagoncu E, Clarke A, Devereux BJ, Tyler LK. Decoding the cortical dynamics of sound-meaning mapping. *J Neurosci*. 2017;37(5):1312–1319. <https://doi.org/10.1523/JNEUROSCI.2858-16.2016>.
- Kujawa SG, Liberman MC. Adding insult to injury: cochlear nerve degeneration after “temporary” noise-induced hearing loss. *J Neurosci*. 2009;29(45):14077–14085. <https://doi.org/10.1523/JNEUROSCI.2845-09.2009>.
- Kumar S, Stephan KE, Warren JD, Friston KJ, Griffiths TD. Hierarchical processing of auditory objects in humans. *PLoS Comput Biol*. 2007;3(6):e100. <https://doi.org/10.1371/journal.pcbi.0030100>.
- Kutas M, Hillyard SA. Event-related brain potentials to grammatical errors and semantic anomalies. *Mem Cogn*. 1983;11(5):539–550. <https://doi.org/10.3758/BF03196991>.
- Ladefoged P. *Vowels and consonants: an introduction to the sounds of languages*. Wiley-Blackwell, Sussex, UK; 2001.
- Ladefoged P, Johnstone K. *A course in phonetics*. Stamford, CT: Cengage Learning; 2015.
- Lee Y-S, Turkeltaub P, Granger R, Raizada RDS. Categorical speech processing in Broca's area: an fMRI study using multivariate pattern-based analysis. *J Neurosci*. 2012;32(11):3942–3948. <https://doi.org/10.1523/JNEUROSCI.3814-11.2012>.
- Leonard MK, Baud MO, Sjerps MJ, Chang EF. Perceptual restoration of masked speech in human cortex. *Nat Commun*. 2016;7(1):13619. <https://doi.org/10.1038/ncomms13619>.
- Lerner Y, Honey CJ, Silbert LJ, Hasson U. Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *J Neurosci*. 2011;31(8):2906–2915. <https://doi.org/10.1523/JNEUROSCI.3684-10.2011>.
- Liberman AM, Cooper FS, Shankweiler DP, Studdert-Kennedy M. Perception of the speech code. *Psychol Rev*. 1967;74(2):431–461. <https://doi.org/10.1037/rev0000013>.
- Liberto GMD, Nie J, Yeaton J, Khalighinejad B, Shamma SA, Mesgarani N. Neural representation of linguistic feature hierarchy reflects second-language proficiency. *NeuroImage*. 2021;227:117586. <https://doi.org/10.1016/j.neuroimage.2020.117586>.
- Lorenzi C, Gilbert G, Carn H, Garnier S, Moore BCJ. Speech perception problems of the hearing impaired reflect inability to use temporal fine structure. *Proc Natl Acad Sci*. 2006;103(49):18866–18869. <https://doi.org/10.1073/pnas.0607364103>.
- Luo H, Poeppel D. Phase patterns of neuronal responses reliably discriminate speech in human auditory cortex. *Neuron*. 2007;54(6):1001–1010. <https://doi.org/10.1016/j.neuron.2007.06.004>.
- Macmillan NA, Kaplan HL. Detection theory analysis of group data: estimating sensitivity from average hit and false-alarm rates. *Psychol Bull*. 1985;98(1):185–199. <https://doi.org/10.1037/0033-2909.98.1.185>.
- Marek S, Tervo-Clemmens B, Calabro FJ, Montez DF, Kay BP, Hatoum AS, Donohue MR, Foran W, Miller RL, Hendrickson TJ, et al. Reproducible brain-wide association studies require thousands of individuals. *Nature*. 2022;603(7902):654–660. <https://doi.org/10.1038/s41586-022-04492-9>.
- Maris E. Enlarging the scope of randomization and permutation tests in neuroimaging and neuroscience. *Biorxiv*. 2019. <https://doi.org/10.1101/685560>.
- Maris E, Oostenveld R. Nonparametric statistical testing of EEG- and MEG-data. *J Neurosci Methods*. 2007;164(1):177–190. <https://doi.org/10.1016/j.jneumeth.2007.03.024>.

- Mattys SL, Jusczyk PW. Phonotactic cues for segmentation of fluent speech by infants. *Cognition*. 2001;78(2):91–121. [https://doi.org/10.1016/S0010-0277\(00\)00109-8](https://doi.org/10.1016/S0010-0277(00)00109-8).
- McDermott JH, Simoncelli EP. Sound texture perception via statistics of the auditory periphery: evidence from sound synthesis. *Neuron*. 2011;71(5):926–940. <https://doi.org/10.1016/j.neuron.2011.06.032>.
- Mesgarani N, Chang EF. Selective cortical representation of attended speaker in multi-talker speech perception. *Nature*. 2012;485(7397):233–236. <https://doi.org/10.1038/nature11020>.
- Mesgarani N, Cheung C, Johnson K, Chang EF. Phonetic feature encoding in human superior temporal gyrus. *Science*. 2014;343(6174):1006–1010. <https://doi.org/10.1126/science.1245994>.
- Millman RE, Johnson SR, Prendergast G. The role of phase-locking to the temporal envelope of speech in auditory perception and speech intelligibility. *J Cogn Neurosci*. 2015;27(3):533–545. https://doi.org/10.1162/jocn_a_00719.
- Moerel M, De Martino F, Santoro R, Ugurbil K, Goebel R, Yacoub E, Formisano E. Processing of natural sounds: characterization of multiplexed spectral tuning in human auditory cortex. *J Neurosci*. 2013;33(29):11888–11898. <https://doi.org/10.1523/JNEUROSCI.5306-12.2013>.
- Moerel M, De Martino F, Kemper VG, Schmitter S, Vu AT, Ugurbil K, Formisano E, Yacoub E. Sensitivity and specificity considerations for fMRI encoding, decoding, and mapping of auditory cortex at ultra-high field. *NeuroImage*. 2018;164:18–31. <https://doi.org/10.1016/j.neuroimage.2017.03.063>.
- Moore BCJ. Perceptual consequences of cochlear hearing loss and their implications for the design of hearing aids. *Ear Hear*. 1996;17(2):133–161. <https://doi.org/10.1097/00003446-199604000-00007>.
- Morosan P, Schleicher A, Amunts K, Zilles K. Multimodal architectonic mapping of human superior temporal gyrus. *Anat Embryol (Berl)*. 2005;210(5-6):401–406. <https://doi.org/10.1007/s00429-005-0029-1>.
- Moulines E, Charpentier F. Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones. *Speech Comm*. 1990;9(5-6):453–467. [https://doi.org/10.1016/0167-6393\(90\)90021-Z](https://doi.org/10.1016/0167-6393(90)90021-Z).
- Narain C, Scott SK, Wise RJS, Rosen S, Leff A, Iversen SD, Matthews PM. Defining a left-lateralized response specific to intelligible speech using fMRI. *Cereb Cortex*. 2003;13(12):1362–1368. <https://doi.org/10.1093/cercor/bhg083>.
- Naselaris T, Kay KN, Nishimoto S, Gallant JL. Encoding and decoding in fMRI. *NeuroImage*. 2011;56(2):400–410. <https://doi.org/10.1016/j.neuroimage.2010.07.073>.
- Naselaris T, Olman CA, Stansbury DE, Ugurbil K, Gallant JL. A voxel-wise encoding model for early visual areas decodes mental images of remembered scenes. *NeuroImage*. 2015;105:215–228. <https://doi.org/10.1016/j.neuroimage.2014.10.018>.
- Norman-Haignere S, Kanwisher Nancy G, McDermott JH. Distinct cortical pathways for music and speech revealed by hypothesis-free voxel decomposition. *Neuron*. 2015;88(6):1281–1296. <https://doi.org/10.1016/j.neuron.2015.11.035>.
- Norman-Haignere SV, Long LK, Devinsky O, Doyle W, Irobunda I, Merricks EM, Feldstein NA, McKhann GM, Schevon CA, Flinker A, et al. Multiscale temporal integration organizes hierarchical computation in human auditory cortex. *Nat Hum Behav*. 2022;6(3):455–469. <https://doi.org/10.1038/s41562-021-01261-y>.
- Nunez-Elizalde AO, Huth AG, Gallant JL. Voxelwise encoding models with non-spherical multivariate normal priors. *NeuroImage*. 2019;197:482–492. <https://doi.org/10.1016/j.neuroimage.2019.04.012>.
- Obleser J, Eisner F, Kotz SA. Bilateral speech comprehension reflects differential sensitivity to spectral and temporal features. *J Neurosci*. 2008;28(32):8116–8123. <https://doi.org/10.1523/JNEUROSCI.1290-08.2008>.
- Oganian Y, Chang EF. A speech envelope landmark for syllable encoding in human superior temporal gyrus. *Sci Adv*. 2019;5(11):eaay6279. <https://doi.org/10.1126/sciadv.aay6279>.
- Overath T, Paik JH. From acoustic to linguistic analysis of temporal speech structure: acousto-linguistic transformation during speech perception using speech quilts. *NeuroImage*. 2021;235:117887. <https://doi.org/10.1016/j.neuroimage.2021.117887>.
- Overath T, Lee JC. The neural processing of phonemes is shaped by linguistic analysis. In: Santurette S, Dau T, Christensen-Dalsgaard J, Tranebjærg L, Andersen T, Poulsen T. (eds.). *Proceedings of the International Symposium on Auditory and Audiological Research*. The Danavox Jubilee Foundation (Ballerup, Denmark). 2017;6:107–116.
- Overath T, McDermott JH, Zarate JM, Poeppel D. The cortical analysis of speech-specific temporal structure revealed by responses to sound quilts. *Nat Neurosci*. 2015;18(6):903–911. <https://doi.org/10.1038/nn.4021>.
- Park H, Ince Robin AA, Schyns Philippe G, Thut G, Gross J. Frontal top-down signals increase coupling of auditory low-frequency oscillations to continuous speech in human listeners. *Curr Biol*. 2015;25(12):1649–1653. <https://doi.org/10.1016/j.cub.2015.04.049>.
- Poeppel D, Idsardi WJ, Vv W. Speech perception at the interface of neurobiology and linguistics. *Philos Trans R Soc B Biol Sci*. 2008;363(1493):1071–1086. <https://doi.org/10.1098/rstb.2007.2160>.
- Preisig BC, Riecke L, Hervais-Adelman A. Speech sound categorization: the contribution of non-auditory and auditory cortical regions. *NeuroImage*. 2022;258:119375. <https://doi.org/10.1016/j.neuroimage.2022.119375>.
- Rao RPN, Ballard DH. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat Neurosci*. 1999;2(1):79–87. <https://doi.org/10.1038/4580>.
- Rauschecker JP, Scott SK. Maps and streams in the auditory cortex: nonhuman primates illuminate human speech processing. *Nat Neurosci*. 2009;12(6):718–724. <https://doi.org/10.1038/nn.2331>.
- Rauschecker JP, Tian B. Processing of band-passed noise in the lateral auditory belt cortex of the rhesus monkey. *J Neurophysiol*. 2004;91(6):2578–2589. <https://doi.org/10.1152/jn.00834.2003>.
- Ringach DL, Sapiro G, Shapley R. A subspace reverse-correlation technique for the study of visual neurons. *Vis Res*. 1997;37(17):2455–2464. [https://doi.org/10.1016/S0042-6989\(96\)00247-7](https://doi.org/10.1016/S0042-6989(96)00247-7).
- Ruggles D, Bharadwaj H, Shinn-Cunningham BG. Normal hearing is not enough to guarantee robust encoding of suprathreshold features important in everyday communication. *Proc Natl Acad Sci USA*. 2011;108(37):15516–15521. <https://doi.org/10.1073/pnas.1108912108>.
- Rutten S, Santoro R, Hervais-Adelman A, Formisano E, Golestani N. Cortical encoding of speech enhances task-relevant acoustic information. *Nat Hum Behav*. 2019;3(9):974–987. <https://doi.org/10.1038/s41562-019-0648-9>.
- Saenz M, Langers DRM. Tonotopic mapping of human auditory cortex. *Hear Res*. 2014;307:42–52. <https://doi.org/10.1016/j.heares.2013.07.016>.
- Saffran JR, Newport EL, Aslin RN. Word segmentation: the role of distributional cues. *J Mem Lang*. 1996;35(4):606–621. <https://doi.org/10.1006/jmla.1996.0032>.

- Samuel AG. Phonemic restoration: insights from a new methodology. *J Exp Psychol Gen.* 1981;110(4):474–494. <https://doi.org/10.1037/0096-3445.110.4.474>.
- Samuel AG. Lexical uniqueness effects on phonemic restoration. *J Mem Lang.* 1987;26(1):36–56. [https://doi.org/10.1016/0749-596X\(87\)90061-1](https://doi.org/10.1016/0749-596X(87)90061-1).
- Santoro R, Moerel M, De Martino F, Goebel R, Ugurbil K, Yacoub E, Formisano E. Encoding of natural sounds at multiple spectral and temporal resolutions in the human auditory cortex. *PLoS Comput Biol.* 2014;10(1):e1003412. <https://doi.org/10.1371/journal.pcbi.1003412>.
- Santoro R, Moerel M, De Martino F, Valente G, Ugurbil K, Yacoub E, Formisano E. Reconstructing the spectrotemporal modulations of real-life sounds from fMRI response patterns. *Proc Natl Acad Sci USA.* 2017;114(18):4799–4804. <https://doi.org/10.1073/pnas.1617622114>.
- Scott SK, Blank CC, Rosen S, Wise RJS. Identification of a pathway for intelligible speech in the left temporal lobe. *Brain.* 2000;123(12):2400–2406. <https://doi.org/10.1093/brain/123.12.2400>.
- Shannon RV, Zeng F-G, Kamath V, Wygonski J, Ekelid M. Speech recognition with primarily temporal cues. *Science.* 1995;270(5234):303–304. <https://doi.org/10.1126/science.270.5234.303>.
- Shin J. Vowels and consonants. In: Brown L, Yeon J. (eds.). *The handbook of Korean linguistics*. Wiley-Blackwell (Sussex, UK); 2015. pp. 3–21.
- Shinn-Cunningham BG, Best V. Selective attention in normal and impaired hearing. *Trends Amplif.* 2008;12(4):283–299. <https://doi.org/10.1177/1084713808325306>.
- Sohn H-M. *The Korean language*. NY: Cambridge University Press; 2001.
- Sohoglu E, Peelle JE, Carlyon RP, Davis MH. Predictive top-down integration of prior knowledge during speech perception. *J Neurosci.* 2012;32(25):8443–8453. <https://doi.org/10.1523/JNEUROSCI.5069-11.2012>.
- Stevens KN. *Acoustic phonetics*. MIT press (Cambridge, MA); 2000.
- Theunissen F, Miller J. Temporal encoding in nervous systems: a rigorous definition. *J Comput Neurosci.* 1995;2(2):149–162. <https://doi.org/10.1007/BF00961885>.
- Vanthonhout J, Decruy L, Wouters J, Simon JZ, Francart T. Speech intelligibility predicted from neural entrainment of the speech envelope. *J Assoc Res Otolaryngol.* 2018;19(2):181–191. <https://doi.org/10.1007/s10162-018-0654-z>.
- Varoquaux G, Raamana PR, Engemann DA, Hoyos-Idrobo A, Schwartz Y, Thirion B. Assessing and tuning brain decoders: cross-validation, caveats, and guidelines. *NeuroImage.* 2017;145(Pt B):166–179. <https://doi.org/10.1016/j.neuroimage.2016.10.038>.
- Verschuere E, Vanthonhout J, Francart T. The effect of stimulus intensity on neural envelope tracking. *Hear Res.* 2021;403:108175. <https://doi.org/10.1016/j.heares.2021.108175>.
- Warren RM. Perceptual restoration of missing speech sounds. *Science.* 1970;167(3917):392–393. <https://doi.org/10.1126/science.167.3917.392>.
- Warren JD, Jennings AR, Griffiths TD. Analysis of the spectral envelope of sounds by the human brain. *NeuroImage.* 2005;24(4):1052–1057. <https://doi.org/10.1016/j.neuroimage.2004.10.031>.
- van de Wiel MA, van Nee MM, Rauschenberger A. Fast cross-validation for multi-penalty high-dimensional ridge regression. *J Comput Graph Stat.* 2021;30(4):835–847. <https://doi.org/10.1080/10618600.2021.1904962>.
- Wild CJ, Davis MH, Johnsrude IS. Human auditory cortex is sensitive to the perceived clarity of speech. *NeuroImage.* 2012;60(2):1490–1502. <https://doi.org/10.1016/j.neuroimage.2012.01.035>.
- Wu MC-K, David SV, Gallant JL. Complete functional characterization of sensory neurons by system identification. *Annu Rev Neurosci.* 2006;29(1):477–505. <https://doi.org/10.1146/annurev.neuro.29.051605.113024>.
- Yi HG, Leonard MK, Chang EF. The encoding of speech sounds in the superior temporal gyrus. *Neuron.* 2019;102(6):1096–1110. <https://doi.org/10.1016/j.neuron.2019.04.023>.
- Yoon T-J, Kang Y. The Korean phonetic aligner program suite. <http://korean.utsc.utoronto.ca/kpa/>. 2013.
- Yuan J, Liberman M. Speaker identification on the SCOTUS corpus. *J Acoust Soc Am.* 2008;123(5_Supplement):3878. <https://doi.org/10.1121/1.2935783>.