

Attentional Modulation of Hierarchical Speech Representations in a Multi-Talker Environment

Ibrahim Kiremitçi^{1,2}, Özgür Yılmaz^{2,3}, Emin Çelik^{1,2}, Mo Shahdloo^{2,4}, Alexander G. Huth^{5,6}, and Tolga Çukur^{1,2,3}

¹Neuroscience Program, Sabuncu Brain Research Center, Bilkent University, Ankara, TR-06800, Turkey

²National Magnetic Resonance Research Center (UMRAM), Bilkent University, Ankara, TR-06800, Turkey

³Department of Electrical and Electronics Engineering, Bilkent University, Ankara, TR-06800, Turkey

⁴Department of Experimental Psychology, Wellcome Centre for Integrative Neuroimaging, University of Oxford, Oxford, OX3 9DU, U.K.

⁵Department of Neuroscience, The University of Texas at Austin, Austin, TX 78712, USA

⁶Department of Computer Science, The University of Texas at Austin, Austin, TX 78712, USA

Running Title: Attentional Modulation of Hierarchical Speech Representations

Correspondence to:

Ibrahim Kiremitçi
National Magnetic Resonance Research Center (UMRAM)
Bilkent University
Ankara, TR-06800, Turkey
TEL: +90 (507) 920-4441
i.kiremitci@bilkent.edu.tr

or

Tolga Çukur
Department of Electrical and Electronics Engineering, Room 304
Bilkent University
Ankara, TR-06800, Turkey
TEL: +90 (312) 290-1164
cukur@ee.bilkent.edu.tr

ABSTRACT

Humans are remarkably adept in listening to a desired speaker in a crowded environment, while filtering out non-target speakers in the background. Attention is key to solving this difficult cocktail-party task, yet a detailed characterization of attentional effects on speech representations is lacking. It remains unclear across what levels of speech features and how much attentional modulation occurs in each brain area during the cocktail-party task. To address these questions, we recorded whole-brain BOLD responses while subjects either passively listened to single-speaker stories, or selectively attended to a male or a female speaker in temporally-overlaid stories in separate experiments. Spectral, articulatory, and semantic models of the natural stories were constructed. Intrinsic selectivity profiles were identified via voxelwise models fit to passive listening responses. Attentional modulations were then quantified based on model predictions for attended and unattended stories in the cocktail-party task. We find that attention causes broad modulations at multiple levels of speech representations while growing stronger towards later stages of processing, and that unattended speech is represented up to the semantic level in parabelt auditory cortex. These results provide insights on attentional mechanisms that underlie the ability to selectively listen to a desired speaker in noisy multi-speaker environments.

Keywords: cocktail-party, dorsal and ventral stream, encoding model, fMRI, natural speech.

INTRODUCTION

Humans are highly adept at perceiving a target speaker in crowded multi-speaker environments (Shinn-Cunningham and Best 2008; Kidd and Colburn 2017; Li et al. 2018). Auditory attention is key to behavioral performance in this difficult “cocktail-party problem” (Cherry 1953; Fritz et al. 2007; McDermott 2009; Bronkhorst 2015; Shinn-Cunningham et al. 2017). Literature consistently reports that attention selectively enhances cortical responses to the target stream in auditory cortex and beyond, while filtering out non-target background streams (Hink and Hillyard 1976; Teder et al. 1993; Alho et al. 1999, 2003, 2014; Jäncke et al. 2001, 2003; Lipschutz et al. 2002; Rienne et al. 2008, 2010; Elhilali et al. 2009; Gutschalk and Dykstra 2014). However, the precise link between the response modulations and underlying speech representations is less clear. Speech representations are hierarchically organized across multiple stages of processing in cortex, with each stage selective for diverse information ranging from low-level acoustic to high-level semantic features (Davis and Johnsrude 2003; Griffiths and Warren 2004; Hickok and Poeppel 2004, 2007; Rauschecker and Scott 2009; Friederici 2011; Di Liberto et al. 2015; de Heer et al. 2017; Brodbeck et al. 2018a). Thus, a principal question is to what extent attention modulates these multi-level speech representations in the human brain during a cocktail-party task (Miller 2016; Simon 2017).

Recent electrophysiology studies on the cocktail-party problem have investigated attentional response modulations for natural speech stimuli (Kerlin et al. 2010; Ding and Simon 2012a, 2012b; Mesgarani and Chang 2012; Power et al. 2012; Zion Golumbic et al. 2013; Puvvada and Simon 2017; Brodbeck et al. 2018b; O’Sullivan et al. 2019; Puschman et al. 2019). Ding and Simon (2012a, 2012b) fit spectrotemporal encoding models to predict cortical responses from the speech spectrogram. Attentional modulation in the peak amplitude of spectrotemporal response functions was reported in planum temporale in favor of the attended speech. Mesgarani and Chang (2012) built decoding models to estimate the speech spectrogram from responses measured during passive listening, and examined the similarity of the decoded spectrogram during a cocktail-party task to the isolated spectrograms of attended versus unattended speech. They found higher similarity to attended speech in non-primary auditory cortex. Zion Golumbic et al. (2013) reported amplitude modulations in speech-envelope response functions towards attended speech across auditory, inferior temporal, frontal and parietal cortices. Other studies using decoding models have similarly reported higher decoding

performance for the speech envelope of the attended stream in auditory, prefrontal, motor and somatosensory cortices (Puvvada and Simon 2017; Puschmann et al. 2019). Brodbeck et al. (2018b) further identified peak amplitude response modulations for sub-lexical features including word onset and cohort entropy in temporal cortex. Note that because these electrophysiology studies fit models for acoustic or sub-lexical features, the reported attentional modulations primarily comprised relatively low-level speech representations.

Several neuroimaging studies have also examined whole-brain cortical responses to natural speech in a cocktail-party setting (Nakai et al. 2005; Alho et al. 2006; Ikeda et al. 2010; Hill and Miller 2010; Wild et al. 2012; Regev et al. 2019). Hill and Miller (2010) measured BOLD response levels while subjects either passively listened to speech streams or attended to a target stream based on pitch or location. Attentional increases in BOLD responses were reported in non-primary auditory cortex as well as insula, frontal and parietal cortices. Furthermore, pitch-based attention was found to elicit higher responses in bilateral posterior and right middle superior temporal sulcus, whereas location-based attention elicited higher responses in left intraparietal sulcus. In alignment with electrophysiology studies, these results suggest that attention modulates relatively low-level speech representations comprising paralinguistic features. In a more recent study, Regev et al. (2019) measured responses under two distinct conditions: while subjects were presented bimodal speech-text stories and asked to attend to either the auditory or visual stimulus, and while subjects were presented unimodal speech or text stories. Correlation of response patterns was measured between unimodal and bimodal conditions. Broad attentional modulations in response correlation were reported from primary auditory cortex to temporal, parietal and frontal regions in favor of the attended modality. While this finding raises the possibility that attention might also affect representations in higher-order regions, a systematic characterization of individual speech features that drive attentional modulations across cortex is lacking.

An equally important question regarding the cocktail-party problem is whether unattended speech streams are represented in cortex despite the reported modulations in favor of the target stream (Bronkhorst 2015; Miller 2016). Electrophysiology studies on this issue identified representations of low-level spectrogram and speech envelope features of unattended speech in early auditory areas (Ding and Simon 2012a, 2012b; Mesgarani and Chang 2012; Zion Golumbic et al. 2013; Puvvada and Simon 2017; Brodbeck et al. 2018b; Puschmann et al. 2019), but no representations of linguistic features (Brodbeck et al. 2018b). Meanwhile, a group of neuroimaging studies found broader cortical responses to unattended speech in superior temporal cortex (Scott et al. 2004, 2009; Wild et al. 2012; Scott and McGettigan 2013; Evans et al. 2016; Regev et al. 2019). Specifically, Wild et al. (2012) and Evans et al. (2016) reported enhanced activity associated with the intelligibility of unattended stream in parts of superior temporal cortex extending to superior temporal sulcus. Although this implies that responses in relatively higher auditory areas carry some information regarding unattended speech stimuli, the specific features of unattended speech that are represented across the cortical hierarchy of speech is lacking.

Here we questioned whether and how attention affects representations of attended and unattended natural speech across cortex. To address these questions, we systematically examined multi-level speech representations during a naturalistic and diotic cocktail-party task. Whole-brain BOLD responses were recorded in two separate experiments while subjects were presented engaging spoken narratives from *The Moth Radio Hour*. In the passive-listening experiment, subjects listened to single-speaker stories for over two hours. Separate voxelwise models were fit that measured selectivity for spectral, articulatory, and semantic features of natural speech during passive listening (de Heer et al. 2017). In the cocktail-party experiment, subjects listened to temporally-overlaid speech streams from two speakers while attending to a target category (male or female speaker). To assess attentional modulation in functional selectivity, voxelwise models fit during passive listening were used to predict

responses for the cocktail-party experiment. Model performances were calculated separately for attended and unattended stories. Attentional modulation was taken as the difference between these two performance measurements. Comprehensive analyses were conducted to examine the intrinsic complexity and attentional modulation of multi-level speech representations and to investigate up to what level of speech features unattended speech is represented across cortex.

Materials and Methods

Participants

Functional data were collected from five healthy adult subjects (four males and one female; aged between 26 and 31) who had no reported hearing problems. The experimental procedures were approved by the Committee for the Protection of Human Subjects at University of California, Berkeley. Written informed consent was obtained from all subjects.

Stimuli

Figure 1 illustrates the two main types of stimuli used in the experiments: single-speaker stories (passive-story) and two-speaker stories (cocktail-story). Ten passive-stories were taken from The Moth Radio Program: “Alternate Ithaca Tom” by Tom Weiser; “How to Draw a Nekkid Man” by Tricia Rose Burt; “Life Flight” by Kimberly Reed; “My Avatar and Me” by Laura Albert; “My First Day at the Yankees” by Matthew McGough; “My Unhurried Legacy” by Kyp Malone; “Naked” by Catherine Burns; “Ode to Stepfather” by Ethan Hawke; “Targeted” by Jen Lee and “Under the Influence” by Jeffery Rudell. All stories were told before a live audience by a male or female speaker, and they were about 10-15 min long. Each cocktail-story was generated by temporally overlaying a pair of stories told by different genders and selected from the passive-story set. When the durations of the two passive stories differed, the longer story was clipped from the end to match durations. Three cocktail-stories were prepared: from “Targeted” and “Ode to Stepfather” (cocktail1); from “How to Draw a Nekkid Man” and “My First Day at the Yankees” (cocktail2); and from “Life Flight” and “Under the Influence” (cocktail3). In the end, the stimuli consisted of ten passive-stories and three cocktail-stories.

Experimental procedures

Figure 1 outlines the two main experiments conducted in separate sessions: passive-listening and cocktail-party experiments. In the passive-listening experiment, subjects were instructed to listen to single-speaker stories vigilantly albeit without any explicit task. Each of the ten passive-stories was presented once in a separate run of the experiment. Two two-hour sessions were conducted, resulting in ten runs of passive-data for each subject. In the cocktail-party experiment, subjects were instructed to listen to two-speaker stories while attending to a target speaker (either the male or the female speaker). Each of the three cocktail-stories was presented twice in separate runs. Different stories were presented in consecutive runs while the attention condition was alternated. An exemplary sequence of runs was: cocktail1-M (attend to male speaker in cocktail1), cocktail2-F (attend to female speaker in cocktail2), cocktail3-M, cocktail1-F, cocktail2-M, and cocktail3-F. The first attention condition assigned was counterbalanced across subjects. A single two-hour session was conducted, resulting in six runs of cocktail-data for each subject. The dataset collected from the passive-listening experiment was previously analyzed (Huth et al. 2016; de Heer et al. 2017); however, the dataset collected from the cocktail-party experiment was specifically collected for this study.

In both experiments, the length of each run was tailored to the length of the story stimulus with additional ten sec of silence both before and after the stimulus. All stimuli were played at 44.1 kHz using Sennheiser S14 in-ear piezo-electric headphones. The frequency response of the headphones was flattened using a Behringer Ultra-Curve Pro Parametric Equalizer.

MRI data collection and preprocessing

MRI data were collected on a 3T Siemens TIM Trio scanner at the Brain Imaging Center, UC Berkeley, using a 32-channel volume coil. For functional scans, a gradient echo EPI sequence was used with TR = 2.0045 s, TE = 31 ms, flip angle = 70°, voxel size = 2.24 x 2.24 x 4.1 mm³, matrix size = 100 x 100, field of view = 224 x 224 mm² and 32 axial slices covering the entire cortex. For anatomical data, a T1-weighted multi-echo MP-RAGE sequence was used with voxel size = 1 x 1 x 1 mm³ and field of view = 256 x 212 x 256 mm³.

Each functional run was motion corrected using FMRIB's Linear Image Registration Tool (FLIRT) (Jenkinson and Smith 2001). A cascaded motion-correction procedure was performed, where separate transformation matrices were estimated within single runs, within single sessions and across sessions sequentially. To do this, volumes in each run were realigned to the mean volume of the run. For each session, the mean volume of each run was then realigned to the mean volume of the first run in the session. Lastly, the mean volume of the first run of each session was realigned to the mean volume of the first run of the first session of the passive-listening experiment. The estimated transformation matrices were concatenated and applied in a single step. Non-brain tissues were removed using Brain Extraction Tool in FSL 5.0 (Smith 2002). Low-frequency drifts in BOLD responses were removed using a 3rd order Savitsky-Golay filter over a 240 s window (Savitzky and Golay 1964). Single voxel responses were then z-scored to have zero mean and unit variance within each run.

Visualization on cortical flatmaps

Cortical flatmaps were used for visualization of model prediction scores, functional selectivity and attentional modulation profiles, and representational complexity and modulation gradients. Cortical surfaces were reconstructed from anatomical data using Freesurfer (Dale et al. 1999). Five relaxation cuts were made into the surface of each hemisphere, and the surface crossing the corpus callosum was removed. Functional data were aligned to the anatomical data via the boundary-based alignment tool in FSL (Greve and Fischl 2009). Voxelwise results were projected onto and visualized on the cortical surface via the pycortex toolbox (Gao et al. 2015).

ROI definitions and abbreviations

We defined region of interests for each subject based on an automatic atlas-based parcellation of the cortex (Destrieux et al. 2010). The cortex was segmented into the regions of the Destrieux atlas (Destrieux et al. 2010) using Freesurfer (Dale et al. 1999); and these anatomical regions were labeled according to the atlas. To explore potential selectivity gradients across the lateral aspects of Superior Temporal Gyrus and Superior Temporal Sulcus, these ROIs were further split into three equidistant sub-regions in posterior-to-anterior direction. We only considered regions with at least ten speech-selective voxels in each individual subject for subsequent analyses.

Table S1 lists the defined ROIs and the number of spectrally, articulatorily and semantically selective voxels within each ROI, with number of speech-selective voxels. ROI abbreviations and corresponding Destrieux indices are: Heschl's Gyrus (HG: 33), Heschl's Sulcus (HS: 74), Planum Temporale (PT: 36), posterior segment of Sylvian Fissure (pSF: 41), lateral aspect of Superior Temporal Gyrus (STG: 34), Superior Temporal Sulcus (STS, 73), Middle Temporal Gyrus (MTG: 38), Angular Gyrus (AG: 25), Supramarginal Gyrus (SMG: 26), Intraparietal Sulcus (IPS: 56), opercular part of Inferior Frontal Gyrus/Pars Opercularis (POP: 12), triangular part of Inferior Frontal Gyrus/Pars Triangularis (PTR: 14), Precentral Gyrus (PreG: 29), medial Occipito-Temporal Sulcus (mOTS:60), Inferior Frontal Sulcus (IFS: 52), Middle Frontal Gyrus (MFG:15), Middle Frontal Sulcus (MFS: 53), Superior Frontal Sulcus (SFS: 54), Superior Frontal Gyrus (SFG: 16), Precuneus (PreC: 30), Subparietal Sulcus (SPS: 71), and Posterior Cingulate Cortex (PCC: 9 and 10). The subregions of STG are: aSTG (anterior one third of STG), mSTG (middle one third of STG) and pSTG (posterior one third of STG). The subregions of STS are: aSTS (anterior one third of STS), mSTS (middle one third of STS) and pSTS (posterior one third of STS). MTG was not split into

subregions since these subregions did not have a sufficient number of speech-selective voxels in each individual subject.

Model construction

To comprehensively assess speech representations, we constructed spectral, articulatory, and semantic models of the speech stimuli (Figure 2; de Heer et al. 2017).

Spectral model. For the spectral model, spectral power density of the sound signal was computed in 300 50-Hz bands between 0 Hz and ~15 kHz. The power spectrum was calculated for 50 ms segments of the sound signal and expressed in dB units. The resulting spectral features were Lanczos-filtered at a cutoff frequency of 0.25 Hz and downsampled to 0.5 Hz to match the sampling rate of fMRI. The 300 spectral features sampled at every two sec were temporally z-scored to zero mean and unit variance.

Articulatory model. For the articulatory model, each phoneme in the stories was mapped onto a unique set of 22 articulation features; for example, phoneme ZH is postalveolar, fricative and voiced (Levelt 1993; de Heer et al. 2017). This mapping resulted in 22-dimensional binary vectors for each phoneme. To obtain the timestamp of each phoneme and word in the stimuli, the speech in the stories were aligned with the story transcriptions using the Penn Phonetics Lab Forced Aligner (Yuan and Liberman 2008). Alignments were manually verified and corrected using Praat (www.praat.org). The articulatory features were Lanczos-filtered at a cutoff frequency of 0.25 Hz and downsampled to 0.5 Hz. Finally, the 22 articulatory features were z-scored to zero mean and unit variance.

Semantic model. For the semantic model, co-occurrence statistics of words were measured via a large corpus of text (Mitchell et al. 2008; Huth et al. 2016; de Heer et al. 2017). The text corpus was compiled from 2,405,569 Wikipedia pages, 36,333,459 user comments scraped from reddit.com, 604 popular books and the transcripts of 13 Moth stories (including the stories used as stimuli). We then built a 10,470-word lexicon from the union set of the 10,000 most common words in the compiled corpus and all words appearing in the ten Moth stories used in the experiment. Basis words were then selected as a set of 985 unique words from Wikipedia's List of 1000 Basic Words. Co-occurrence statistics of the lexicon words with 985 basis words within a 15-word window were characterized as a co-occurrence matrix of size 985x10,470. Elements of the resulting co-occurrence matrix were log-transformed, z-scored across columns to correct for differences in basis-word frequency, and z-scored across rows to correct for differences in lexicon-word frequency. Each word in the stimuli was then represented with a 985-dimensional co-occurrence vector based on the speech-transcription alignments. The semantic features were Lanczos-filtered at a cutoff frequency of 0.25 Hz and downsampled to 0.5 Hz. The 985 semantic features were finally z-scored to zero mean and unit variance.

Decorrelation of feature spaces. In natural stories, there might be potential correlations among certain spectral, articulatory, or semantic features. If significant, such correlations can partly confound assessments of model performance. To assess the unique contribution of each feature space to the explained variance in BOLD responses, a decorrelation procedure was first performed (Figure 2). To decorrelate a feature matrix F of size $m \times n$ from a second feature matrix K of size $m \times p$, we first found an orthonormal basis for the column space of K ($col\{K\}$) using economy-size singular value decomposition:

$$K_{m \times p} = U_{m \times p} \times S_{p \times p} \times V_{p \times p}$$

where U contains left singular vectors as columns, V contains right singular vectors, and S contains the singular values. Left singular vectors were taken as the orthonormal basis for $col\{K\}$, and each column of F was decorrelated from it according to the following formula:

$$\vec{f}_i^d = \vec{f}_i - \sum_{j=1}^p (\vec{f}_i \cdot \vec{u}_j) \cdot \vec{u}_j$$

Where \vec{f}_i , \vec{u}_j are the column vectors of F and U respectively, and \vec{f}_i^d is the column vectors of the decorrelated feature matrix, F^d . To decorrelate feature matrices for the models considered here, we took the original articulatory feature matrix as a reference, and decorrelated the spectral feature matrix from the articulatory feature matrix, and decorrelated the semantic feature matrix from both articulatory and spectral feature matrices. This decorrelation sequence was selected because spectral and articulatory features capture lower-level speech representations, and the articulatory feature matrix had the fewest number of features among all models. In the end, we obtained 3 decorrelated feature matrices whose columns had zero correlation with the columns of the other two matrices.

Analyses

The main motivation of this study is to understand whether and how strongly various levels of speech representations are modulated across cortex during a cocktail-party task. To answer this question, we followed a two-stage approach as illustrated in Figure 3. In the first stage, we identified voxels selective for speech features using data from the passive-listening experiment. To do this, we measured voxelwise selectivity separately for spectral, articulatory, and semantic features of the passive-stories. In the second stage, we used the models fit using passive-data to predict BOLD responses measured in the cocktail-party experiment. Prediction scores for attended versus unattended stories were compared to quantify the degree of attentional modulations, separately for each model and globally across all models.

Note that a subset of the ten passive-stories were used to generate three cocktail-stories used in the experiments. Hence, to prevent potential bias, a three-fold cross-validation procedure was performed for testing models fit using passive-data on cocktail-data. In each fold, models were fit using eight-run passive-data; and separately tested on two-run passive-data and two-run cocktail-data. There was no overlap between the stories in the training and testing runs. Model predictions were aggregated across three folds, and prediction scores were then computed.

Voxelwise modeling. In the first stage, we fit voxelwise models in individual subjects using passive-data. To account for hemodynamic delays, we used a linearized four-tap finite impulse response (FIR) filter to allow different HRF shapes for separate brain regions (Goutte et al. 2000). Each model feature was represented as four features in the stimulus matrix to account for their delayed effects in BOLD responses at 2, 4, 6 and 8 sec. Model weights, W , were then found using L2-regularized linear regression:

$$W = (F^T F + I\lambda)^{-1} F^T R$$

Here, λ is the regularization parameter, F is the decorrelated feature matrix for a given model and R is the aggregate BOLD response matrix for cortical voxels. A cross-validation procedure with 50 iterations was performed to find the best regularization parameter for each voxel among 30 equispaced values in log-space of $1:10^5$. The training passive-data was split into 50 equisized chunks, where 1 chunk was reserved for validation and 49 chunks were reserved for model fitting at each iteration. Prediction scores were taken as Pearson's correlation between predicted and measured BOLD responses. The optimal λ value for each voxel was selected by maximizing the average prediction score across cross-validation folds. The final model weights were obtained using the entire set of training passive-data and the optimal λ . Next, we measured the prediction scores of the fit models on testing passive-data. When a model had less than ten significantly predicted voxels within an ROI, those voxels were excluded for that ROI and that model. Speech-selective voxels were then taken as the union of

voxels significantly predicted by spectral, articulatory, and semantic models ($q(FDR) < 10^{-4}$; t -test). Subsequent analyses were performed on the speech-selective voxels.

a) Model-specific selectivity index. Single-voxel prediction scores on passive-data were used to quantify the degree of selectivity of each ROI to the underlying model features under passive-listening. To do this, a model-specific selectivity index, (SI_m), was defined as follows:

$$SI_m = \frac{(r)_m}{\sum_i (r)_i}, \quad i, m \in \{spe, art, sem\}$$

where r is the average prediction score across speech-selective voxels within the ROI during passive-listening. SI_m is in the range of [0, 1], where higher values indicate stronger selectivity for the underlying model.

b) Complexity index. The complexity of speech representations was characterized via a complexity index, (CI), which reflected the relative tuning of an ROI for low- versus high-level speech features. The following intrinsic complexity levels were assumed for the three speech models considered here: ($c_{spe}, c_{art}, c_{sem}$) = (0.0, 0.5, 1.0). Afterwards, CI was taken as the average of the complexity levels weighted by the selectivity indices:

$$CI = \sum_m SI_m c_m, \quad m \in \{spe, art, sem\}$$

CI is in the range of [0, 1], where higher values indicate stronger tuning for semantic features and lower values indicate stronger tuning for spectral features.

Assessment of attentional modulations. In the second stage, we tested the passive-models on cocktail-data to quantify ROI-wise attentional modulation in selectivity for corresponding model features and to find the extent of the representation of unattended speech. These analyses were repeated separately for the three speech models.

a) Model-specific attention index. To quantify the attentional modulation in selectivity for speech features, we compared prediction scores for attended versus unattended stories in the cocktail-party experiment. Models fit using passive-data were used to predict BOLD responses elicited by cocktail-stories. In each run, only one of the two speakers in a cocktail-story was attended while the other speaker was ignored. Separate response predictions were obtained using the isolated story stimuli for the attended and unattended speakers. Since a voxel can represent information on both attended and unattended stimuli, a weighted linear combination of these predicted responses was considered:

$$R_c = R_a w_c + R_u (1 - w_c),$$

Where R_a and R_u are the predicted responses for the attended and unattended stories in a given run; R_c is the combined response and w_c is the combination weight. We computed R_c for each separate w_c value in [0:0.1:1]. Note that $R_c = R_a$ when $w_c = 1.0$; and $R_c = R_u$ when $w_c = 0.0$. We then calculated single-voxel prediction scores for each w_c value. An illustrative plot of r_c/r_{max} vs w_c is given in Figure 3b, where r_c denotes the prediction scores and r_{max} denotes the maximum r_c value (the optimal combination). r_a and r_u are the prediction scores for attended and unattended stories respectively. To quantify the degree of attentional modulation, a model-specific attention index (AI_m) was taken as:

$$AI_m = \alpha_m \left(\frac{r_a - r_u}{r_{max}} \right)_m, \quad \alpha_m = \frac{(r_{max})_m}{\sum_i (r_{max})_i}, \quad m, i \in \{spe, art, sem\},$$

where r_{max} denotes an ideal upper limit for model performance, and α_m reflects the relative model performance under the cocktail-party task. Note that AI_m considers selectivity to the underlying model features when calculating the degree of attentional modulation.

b) Global attention index. We then computed global attention index (gAI) as follows:

$$gAI = \sum_m AI_m, \quad m \in \{spe, art, sem\}$$

Both gAI and AI_m are in the range $[-1,1]$. A positive index indicates attentional modulation of selectivity in favor of the attended stimuli and a negative index indicates attentional modulation in favor of the unattended stimuli. A value of zero indicates no modulation.

Statistical Tests

Significance assessments within subjects. For each voxel-wise model, significance of prediction scores was assessed via a t-test; and resulting p-values were false-discovery-rate corrected for multiple comparisons (FDR; Benjamini and Hochberg 1995).

A bootstrap test was used in assessments of SI_m , CI , AI_m and gAI within single subjects. In ROI analyses, speech-selective voxels within a given ROI were resampled with replacement 10000 times. For each bootstrap sample, mean prediction score of a given model was computed across resampled voxels. Significance level was taken as the fraction of bootstrap samples in which the test metric computed from these prediction scores is less than 0 (for right-sided tests) or greater than 0 (for left-sided tests). The same procedure was also used for comparing pairs of ROIs, where ROI voxels were resampled independently.

Significance assessments across subjects. A bootstrap test was used in assessments of SI_m , CI , AI_m and gAI across subjects. In ROI analyses, ROI-wise metrics were resampled across subjects with replacement 10000 times. Significance level was taken as the fraction of bootstrap samples where the test metric averaged across resampled subjects is less than 0 (for right-sided tests) or greater than 0 (for left-sided tests). The same procedure was also used for comparisons among pairs of ROIs.

Results

Attentional modulation of multi-level speech representations

Recent electrophysiology (Ding and Simon 2012a, 2012b; Power et al. 2012; Zion Golumbic et al. 2013; Brodbeck et al. 2018b; O'Sullivan et al. 2019) and neuroimaging studies (Hill and Miller 2010; Wild et al. 2012; Regev et al. 2019) suggest that attention modulates cortical responses to speech during non-spatial or spatial cocktail-party tasks. However, less is known regarding the cortical distribution and strength of these modulations for features involved in speech perception. To examine this issue, we first obtained a baseline measure of intrinsic selectivity for speech features. For this purpose, we fit voxelwise models using BOLD responses recorded during passive listening. Speech representations are thought to be organized hierarchically across multiple stages of processing in the brain, ranging from acoustic features in early auditory cortex to linguistic features in downstream areas (Davis and Johnsrude 2003; Hickok and Poeppel 2007; Okada et al. 2010; DeWitt and Rauschecker 2012; Bizley and Cohen 2013). To broadly examine this hierarchy, we built three separate models containing low-

level spectral, intermediate-level articulatory and high-level semantic features of natural stories (de Heer et al. 2017). Supplementary Fig. S1 displays the cortical distribution of prediction scores for each model in a representative subject, and Supplementary Table S1 lists the number of significantly predicted voxels by each model in anatomical ROIs. We find *spectrally-selective voxels* mainly in early auditory regions (bilateral HG, HS and PT; and left pSF) and bilateral SMG, and *articulatorily-selective voxels* mainly in early auditory regions (bilateral HG, HS and PT; and left pSF), bilateral STG, STS, SMG and MFS as well as left POP and PreG. In contrast, *semantically-selective voxels* are found broadly across cortex except early auditory regions (bilateral HG and HS; and right PT).

To quantitatively examine overlap among spectral, articulatory, and semantic representations in cortex, we separately measured the degree of functional selectivity for each feature level via a model-specific selectivity index (SI_m ; see Methods). Bar plots of selectivity indices are displayed in Supplementary Fig. S2a for perisylvian cortex and in Supplementary Fig. S3 for non-perisylvian cortex. Several distinct selectivity profiles are observed from distributed selectivity for spectral, articulatory, and semantic features (e.g., right SMG) to strong tuning to a single level of features (e.g., left IPS and AG). Cortical ROIs were separated into five characteristic groups based on their selectivity profiles (Supplementary Fig. S2b). In Profile-1 (P1), both SI_{art} and SI_{spe} are dominant ($p < 10^{-3}$); in P2 all indices are dominant ($p < 10^{-4}$); in P3, SI_{art} is dominant ($p < 10^{-2}$); in P4, SI_{art} and SI_{sem} are dominant ($p < 0.05$); and in P5, SI_{sem} is dominant ($p < 0.05$) (see Supplementary Table S2a for detailed significance tests). A progression from P1 to P4 is apparent while moving from primary auditory cortex to intermediate regions in temporal cortex, with P5 primarily manifesting in higher-order regions. To examine the hierarchical organization of the speech representations in a finer scale, we also defined a complexity index, CI , that reflects whether an ROI is relatively tuned for low-level spectral or high-level semantic features. A detailed investigation of the gradients in CI across two main auditory streams (dorsal and ventral stream) was conducted (see Supplementary Results). These results corroborate the view that speech representations are hierarchically organized across cortex with partial overlap mostly in early and intermediate stages of speech processing.

Next, we systematically examined attentional modulations at each level of speech representation during a diotic cocktail-party task. To do this, we recorded whole-brain BOLD responses while participants listened to temporally-overlaid spoken narratives from two different speakers and attended to either a male or female speaker in these two-speaker stories. We used the spectral, articulatory, and semantic models fit using passive-data to predict responses during the cocktail-party task. Since a voxel can represent information on both attended and unattended stimuli, response predictions were expressed as a convex combination of individual predictions for the attended and unattended story within each cocktail-story. Prediction scores were computed based on estimated responses as the combination weights were varied in [0 1] (see Methods). Scores for the optimal combination model were compared against the scores from the individual models for attended and unattended stories. If the optimal combination model significantly outperforms the individual models, it indicates that the voxel represents information from both attended and unattended stimuli. In contrast, if the optimal combination model performs similarly to the individual model for the attended story, the voxel does not represent significant information on the unattended story.

Figure 4 displays prediction scores of the spectral, articulatory, and semantic models as a function of the combination weight in representative ROIs. Scores based on only attended story (r_a), based on only the unattended story (r_u), and based on the optimal combination of the two (r_{max}) are marked. A diverse set of attentional effects are apparent in HS, HG and PT. For the *spectral model* in left HS, the optimal combination puts matched weights to attended and unattended stories, while no significant difference exists between r_a and r_u ($p > 0.05$). This finding implies that attention does not influence spectral representations in left HS. For the *articulatory model* in left HG, r_a is larger than r_u ($p < 10^{-4}$), while r_{max} is greater than r_a ($p < 10^{-2}$). This result suggests that attention impacts articulatory

representations mildly in left HG such that articulatory representations of the unattended story are still maintained to an extent. For the *semantic model* in left PT, r_a is greater than r_u ($p < 10^{-4}$), while no significant difference exists between r_{max} and r_a ($p > 0.05$). This finding indicates that attention affects semantic representations in left PT strongly such that no trace of semantic representations of unattended story is found. A simple inspection of these results suggests that attention may have distinct effects at various levels of speech representation across cortex. Hence, a detailed quantitative analysis is warranted to measure the effect of attention at each level.

Level-specific attentional modulations. To quantitatively assess the strength and direction of attentional modulations, we separately investigated the modulatory effects on spectral, articulatory, and semantic features across cortex. To measure modulatory effects at each feature level, a model-specific attention index (AI_m) was computed, reflecting the difference in model prediction scores when the stories were attended versus unattended (see Methods). AI_m is in the range of [-1, 1]; a positive index indicates selectivity modulation in favor of the attended stimulus, whereas a negative index indicates selectivity modulation in favor of the unattended stimulus. A value of zero indicates no modulation.

Figure 5a and Supplementary Fig. S6 displays the attention index for the spectral, articulatory, and semantic models across perisylvian and non-perisylvian ROIs, respectively. Here we discuss the attention index for each model individually. *Spectral modulation* starts in HG ($p < 0.01$) bilaterally and extends to HS only in the right hemisphere (RH; $p < 0.01$). *Articulatory modulation* also starts as early as HG bilaterally ($p < 10^{-4}$). In the dorsal stream, it extends to PreG and POP in the left hemisphere (LH) and to SMG in the right hemisphere (RH; $p < 10^{-4}$). In the ventral stream, it extends to PTR and MTG bilaterally ($p < 10^{-3}$). Articulatory modulation is also apparent -albeit generally less strong- in some frontal regions (bilateral IFS, MFG, MFS, SFG) and parietal regions (bilateral SPS, right PrC and SPS) ($p < 0.05$). In the dorsal stream, *semantic modulation* starts in PT and extends to POP and PreG in LH, whereas it only occurs in SMG in RH ($p < 10^{-3}$). In the ventral stream, it extends from mSTG to MTG and PTR bilaterally ($p < 0.01$). Lastly, strong semantic modulation is observed widespread across higher-order regions within frontal, parietal, and occipital cortices ($p < 0.01$). Taken together, these results suggest that attending to a target speaker causes broad selectivity modulations distributed across cortex at the linguistic level (articulatory and semantic), yet modulations at the acoustic level (spectral) are primarily constrained to early auditory cortex.

Attention is postulated to be a multi-level selection process that affects each stage of stimulus processing disparately based on its intrinsic function (Kastner and Pinsk 2004). Previous studies on visual attention have reported that attention predominantly enhances representations in areas that are preferentially selective to target features (Maunsell and Treue 2006; Carrasco 2011). Thus, we hypothesized that attending to a target speaker in a cocktail-story would manifest stronger modulations in a region for features that the ROI was preferentially selective for during passive listening. To test this hypothesis, we examined the dominant attention index across the three models within each ROI as shown in Figure 5b. Among ROIs with intrinsic selectivity profile P1 (dominant SI_{spe} and SI_{art}), AI_{art} and AI_{spe} are dominant in bilateral HG and right HS, whereas only AI_{art} is dominant in left HS and pSF, and right PT ($p < 0.01$). For P2 (dominant SI_{spe} , SI_{art} and SI_{sem}) and P4 (dominant SI_{art} and SI_{sem}), AI_{art} and AI_{sem} are dominant ($p < 0.01$; but see left aSTG where only AI_{sem} is dominant). For P3 (dominant SI_{art}), AI_{art} is dominant ($p < 0.05$). For P5 (dominant SI_{sem}), AI_{sem} is dominant in all ROIs ($p < 0.05$) except for right IFS and left POP where both AI_{sem} and AI_{art} are dominant ($p < 10^{-4}$).

To quantify the degree of similarity between the intrinsic selectivity and attentional modulation profiles, we further measured the cosine similarity between the selectivity and modulation vectors across the three models, $\vec{SI} = (SI_{spe}, SI_{art}, SI_{sem})$ and $\vec{AI} = (AI_{spe}, AI_{art}, AI_{sem})$ (see Supplementary Fig. S7). In all regions, the two vectors show strong similarity ($r > 0.81$, $p < 10^{-4}$, bootstrap test). Correlation between the selectivity and modulation vectors is less strong in several regions including left HS and

right PT. The relatively lower correlation can be attributed to the lack of spectral modulation despite the intrinsic tuning for spectral features. Therefore, our results support the view that attention enhances speech representations to favor target features that each brain area is intrinsically selective for during passive listening, albeit this effect is relatively weaker at the spectral level.

Global attentional modulations. It is commonly assumed that attentional effects grow stronger towards higher-order regions across the cortical hierarchy of speech (Zion Golumbic et al. 2013; O'Sullivan et al. 2019; Regev et al. 2019). Yet, a systematic examination of attentional modulation gradients across dorsal and ventral streams is lacking. To examine this issue, we measured overall attentional modulation in each region via a global attention index (gAI ; see Methods). Similar to the model-specific attention indices, a positive gAI indicates modulations in favor of the attended stimulus, and a negative gAI indicates modulations in favor of the unattended stimulus.

a) Dorsal stream. We first examined variation of gAI across the dorsal stream (left dorsal-1: $HG_L \rightarrow HS_L \rightarrow PT_L \rightarrow (SMG_L) \rightarrow POP_L$, left dorsal-2: $HG_L \rightarrow HS_L \rightarrow PT_L \rightarrow (SMG_L) \rightarrow PreG_L$, and right dorsal: $HG_R \rightarrow HS_R \rightarrow PT_R \rightarrow SMG_R$) as shown in Figure 6. We find significant increase in gAI across the following subtrajectories ($p < 0.01$): $gAI_{HS} < gAI_{PT} < gAI_{SMG} < gAI_{POP}$ and $gAI_{HS} < gAI_{PT} < gAI_{SMG} < gAI_{PreG}$ in the left dorsal stream, and $gAI_{PT} < gAI_{SMG}$ in the right dorsal stream. In contrast, we find no difference between HG and HS in the left dorsal stream ($p > 0.05$), and a significant decrease in gAI from HG to PT in the right dorsal stream ($p < 0.05$). These results suggest that attentional modulations grow progressively stronger across the dorsal stream in LH, whereas the modulation patterns are less consistent in the RH.

b) Ventral stream. We then examined variation of gAI across the ventral stream (left ventral-1: $HG_L \rightarrow mSTG_L \rightarrow mSTS_L \rightarrow MTG_L$, left ventral-2: $HG_L \rightarrow mSTG_L \rightarrow aSTG_L \rightarrow PTR_L$, right ventral-1: $HG_R \rightarrow mSTG_R \rightarrow mSTS_R \rightarrow MTG_R$ and right ventral-2: $HG_R \rightarrow mSTG_R \rightarrow aSTG_R \rightarrow PTR_R$), as shown in Figure 6. We find significant increase in gAI across the following subtrajectories ($p < 10^{-4}$): $gAI_{HG} < gAI_{mSTG} < gAI_{aSTG}$ and $gAI_{HG} < gAI_{mSTG} < gAI_{mSTS}$ in the left ventral stream, and $gAI_{mSTG} < gAI_{aSTG}$ in the right ventral stream. In contrast, we find no difference between aSTG and PTR bilaterally, between mSTS and MTG in the left ventral stream, and between HG, mSTG, mSTS and MTG in the right ventral stream ($p > 0.05$). These results suggest that overall attentional modulations gradually increase across the ventral stream, and that the increases are more consistent in LH compared to RH.

c) Representational complexity versus attentional modulation. Visual inspection of Supplementary Fig. S5b and Figure 6b suggests that the subtrajectories with significant increases in CI and in gAI largely overlap in LH. To quantitatively examine this overlap, we analyzed the correlation between CI and gAI across the subtrajectories where significant increases in CI are obtained. We find significant correlations in $HS \rightarrow PT \rightarrow PreG$ ($r = 0.87$, $p < 10^{-4}$, bootstrap test), in $HS \rightarrow PT \rightarrow SMG \rightarrow POP$ ($r = 0.84$, $p < 10^{-4}$), in $HG \rightarrow mSTG \rightarrow mSTS \rightarrow MTG$ ($r = 0.86$, $p < 10^{-4}$) and in $HG \rightarrow mSTG \rightarrow aSTG$ ($r = 0.99$, $p < 10^{-4}$). In line with a recent study arguing for stronger attentional modulation and higher representational complexity in STG compared to HG (O'Sullivan et al. 2019), our results indicate that attentional modulation increases towards higher-order regions as the representational complexity increases across the dorsal and ventral streams in LH.

d) Hemispheric asymmetries in attentional modulation. To assess potential hemispheric asymmetries in attentional modulation, we compared gAI between the left and right counterparts of each ROI. This analysis was restricted to ROIs with consistent selectivity for speech features in both hemispheres in each individual subject (see Methods). Supplementary Table S3 lists the results of this across-hemisphere comparison. Among ROIs with selectivity profiles P1-P4, higher gAI in LH is observed for PT, aSTG, mSTG, aSTS, mSTS, pSTS, whereas higher gAI in RH is apparent for HG ($p < 0.05$). For selectivity profile P5, higher gAI in LH is observed for PTR, IFS and MFG, whereas higher gAI in

RH is observed for IPS ($p < 0.05$). These results indicate that attentional modulations are right lateralized in earlier stages of speech processing, whereas modulations are left lateralized in intermediate and higher stages (Xiang et al. 2010; Power et al. 2012; Brodbeck et al. 2018b); albeit lateralization is weaker towards later stages.

Cortical representation of unattended speech

An important question regarding multi-speaker speech perception is to what extent unattended stimuli are represented in cortex. Several neuroimaging studies on this topic have reported responses to the unattended speech stream in superior temporal cortex (Scott et al. 2004, 2009; Wild et al. 2012; Evans et al. 2016; Regev et al. 2019). While Wild et al. (2012) and Evans et al. (2016) have further suggested that the responses increase with the intelligibility of the unattended stream in parts of this area, the specific features and brain areas that mediate representation of unattended speech remain unclear. To address this question, here we investigated spectral, articulatory, and semantic representations of unattended stories during the cocktail-party task. We reasoned that if significant information about unattended speech is represented in a brain region, then features of unattended speech should explain significant variance in measured BOLD responses. To test this, we compared the prediction score of a combination model comprising the features of both attended and unattended stories (optimal convex combination) against the prediction score of an individual model comprising only the features of the attended story (see Methods). If the combination model significantly outperforms the individual model in an ROI, then the corresponding features of unattended speech are significantly represented in that ROI.

Figure 7 displays prediction scores based on the features of only the attended story and the optimal convex combination of the attended and unattended stories for each ROI in the dorsal and ventral stream in RH (see Supplementary Fig. S9 for LH). Along the left (HG → HS → PT → SMG → (POP, PreG)) and right (HG → HS → PT → SMG) dorsal stream, spectral features of unattended speech are represented bilaterally up to SMG ($p < 0.01$), articulatory features are represented up to PT in LH and SMG in RH ($p < 0.05$), whereas no semantic representation is apparent ($p > 0.05$). Along the left ventral stream (HG → mSTG → mSTS → MTG and HG → mSTG → aSTG → PTR), spectral features are represented in HG ($p < 10^{-4}$), and articulatory features are represented up to mSTG ($p < 0.05$), again with no semantic representation ($p > 0.05$). In the right ventral stream (HG → mSTG → mSTS → MTG and HG → mSTG → aSTG → PTR), spectral features are represented up to aSTG and mSTS ($p < 0.05$), articulatory features are represented up to mSTG ($p < 0.01$), and semantic representations are found in mSTS ($p < 0.05$). These results indicate that cortical representations of unattended speech in multi-speaker environments extend from the spectral to the semantic level, albeit semantic representations are constrained to right parabelt auditory cortex (mSTS). Furthermore, representations of unattended speech are more broadly spread across the right hemisphere. Note that prior studies have reported response correlations and anatomical overlap between these belt/parabelt auditory regions and the reorienting attention system in the right-hemisphere (Corbetta et al. 2008; Vossel et al. 2014; Puschmann et al. 2017). Therefore, relatively broader representations of unattended speech in the right hemisphere might facilitate distractor detection and filtering during auditory attention tasks.

Discussion

In this study, we investigated the effects of auditory attention on multi-level speech representations across cortex during a diotic and naturalistic cocktail-party task. To assess baseline selectivity for multi-level speech features, we first fit spectral, articulatory, and semantic models using responses recorded during passive listening. We then quantified the complexity of intrinsic representations in each brain region. Next, we used fit models that reflect baseline selectivity for speech features to assess

attentional modulation of speech representations. To do this, responses predicted using stimulus features of attended and unattended stories were compared with responses recorded during the cocktail-party task. This study is among the first to quantitatively characterize attentional modulations in multi-level speech representations of attended and unattended stimuli across speech-related cortex.

Attentional Modulations

The effects of auditory attention on cortical responses have been primarily examined in the literature using controlled stimuli such as simple tones, melodies and isolated syllables or words (Alho et al. 1999; Jäncke et al. 2001, 2003; Lipschutz et al. 2002; Petkov et al. 2004; Johnson and Zatorre 2005; Degerman et al. 2006; Rinne et al. 2005, 2008, 2010; Woods et al. 2009, 2010; Paltoglou et al. 2009; Da costa et al. 2013; Seydell et al. 2014; Riecke et al. 2017). As such, less is known regarding how attention alters hierarchical representations of natural speech. Recent studies on this topic have reported attentional modulations of low-level speech representations comprising acoustic and sub-lexical features during the cocktail-party task (Ding and Simon 2012a, 2012b; Mesgarani and Chang 2012; Zion Golumbic et al. 2013; Puvvada and Simon 2017; Brodbeck et al. 2018b; Puschmann et al. 2019). Going beyond prior reports, here we find that attentional modulations are not solely constrained to the acoustic and sub-lexical levels but also extend to the higher semantic level. Importantly, attending to a male or female speaker categorically causes modulations at all examined feature levels embodied by the target stimulus. This finding is consistent with visual attention studies suggesting that category-based attention influences selectivity for all subordinate features of a target visual object (O'Craven et al. 1999; Shinn-Cunningham 2008).

Several prior studies have reported attentional response modulations for speech-envelope and spectrogram features in non-primary auditory cortex and even higher-order areas (Ding and Simon 2012a, 2012b; Mesgarani and Chang 2012; Zion Golumbic et al. 2013). While our results indicate that attentional modulations for articulatory and semantic representations distribute broadly across cortex, we find that modulations for spectral representations are mainly constrained to the primary auditory cortex. Note that speech envelope and spectrogram features in natural speech carry intrinsic information about linguistic features including syllabic boundaries and articulatory features (Ding and Simon 2014; Liberto et al. 2015). These stimulus correlations can render it challenging to dissociate unique selectivity for articulatory versus spectral features. To minimize biases from potential stimulus correlations, here we leveraged a decorrelation procedure to obtain orthogonal spectral, articulatory, and semantic feature matrices for the stimulus. The distinct modeling procedure for natural speech features might have contributed to the disparities between the current and previous studies on the cortical extent of spectral modulations.

Attention is taken to be a multi-level selection process that affects each stage of stimulus processing uniquely based on its intrinsic functional selectivity (Kastner and Pinsk 2004). Here we first measured the selectivity profiles of brain areas across three distinct feature spaces, from spectral to semantic. We also measured attentional modulation profiles across the same set of feature spaces. A comparison of the selectivity and attentional modulation profiles reveals that attention enhances speech representations to predominantly favor features that each brain area is intrinsically selective for. This finding is consistent with recent studies on auditory attention reporting in primary auditory cortex that the strongest response modulations occur in populations that are preferentially selective to the target feature (Paltoglou et al. 2009; Da Costa et al. 2013; Riecke et al. 2017).

An important question regarding auditory attention is how the strength of attentional affects are distributed across cortex. A common view is that attentional modulations grows relatively stronger towards later stages of processing (Golumbic et al. 2013). Recent studies support this view by reporting bilaterally stronger modulations in frontal versus temporal cortex (Regev et al. 2019) and in non-primary versus primary auditory cortex (O'Sullivan et al. 2019). Adding to this body of evidence, we

further show that attentional modulations gradually increase across the dorsal and ventral streams in the left hemisphere, as the complexity of speech representations grow. While a similar trend is observed across the right hemisphere, several exceptions are reported in belt and parabelt auditory regions including PT where the gradient in attentional modulation is less consistent. Furthermore, attentional modulations are weaker in the right versus left hemisphere across these auditory regions. Note that belt and parabelt regions are suggested to be connected to the right temporo-parietal junction (TPJ) during selective listening (Puschmann et al. 2017). TPJ is one of the central nodes in the reorienting attention system that monitors salient events to filter out distractors and help maintaining focused attention (Corbetta and Schulman 2002; Corbetta et al. 2008; Vossel et al. 2014). Hence less consistent gradients and relatively weaker attentional modulations in belt and parabelt auditory regions in the right hemisphere might suggest a functional role in detecting salient events in the unattended stream during selective listening tasks.

Representation of the unattended speech

Whether unattended speech is represented in cortex during selective listening and if so, at what feature levels its representations are maintained are crucial aspects of auditory attention. Behavioral accounts suggest that unattended speech is primarily represented at the acoustic level (Cherry 1953; Broadbent 1958). Corroborating these accounts, recent electrophysiology studies have identified acoustic representations of unattended speech localized to auditory cortex (Ding and Simon 2012a, 2012b; Zion Golumbic et al. 2013; Puvvada and Simon et al. 2017; Brodbeck et al. 2018b; O'Sullivan et al. 2019; Puschman et al. 2019). In contrast, here we find that acoustic representations of unattended speech extend beyond the auditory cortex as far as SMG in the dorsal stream. Because SMG partly overlaps with the reorienting attention system, unattended speech representations in this region might contribute to filtering of distractors during the cocktail-party task (Corbetta et al. 2008; Vossel et al. 2014).

A more controversial discussion is focused on whether unattended speech representations carry information at the linguistic level (Driver 2001; Lavie 2005; Boulenger et al. 2010; Bronkhorst 2015; Kidd and Colburn 2017). Prior studies on this issue are split between those suggesting the presence (Wild et al. 2012; Evans et al. 2016) versus absence (Sabri et al. 2008; Brodbeck et al. 2018b) of linguistic representations. Here, we find that articulatory representations of unattended speech extend to belt/parabelt auditory areas in the left dorsal and right ventral stream, and to inferior parietal cortex in the right dorsal stream. We further find semantic representation of unattended speech in the right ventral stream (mSTS). These linguistic representations of unattended speech are naturally weaker than those of attended speech, and they are localized to early-to-intermediate stages of auditory processing. Our findings suggest that unattended speech is represented at the linguistic level prior to entering the broad semantic system where full selection of the attended stream occurs (Bregman 1994; Pulvermüller and Shtyrov 2006; Relander et al. 2009; Näätänen et al. 2011; Rämä et al. 2012; Bronkhorst 2015; Ding et al. 2018). Overall, these linguistic representations might serve to direct exogenous triggering of attention to salient features in unattended speech (Moray 1959; Treisman 1960, 1964; Wood and Cowan 1995; Driver 2001; Bronkhorst 2015). Meanwhile, attenuated semantic representations in the ventral stream might facilitate semantic priming of the attended stream by relevant information in the unattended stream (Lewis 1970; Driver 2001; Rivenez et al. 2006).

Conclusion

In sum, our results indicate that attention during a naturalistic cocktail-party task gradually selects attended over unattended speech across both dorsal and ventral processing pathways. This selection is mediated by representational modulations extending from acoustic to linguistic features, where strongest modulations at each stage of processing occur for features that the stage is intrinsically selective for. Despite broad attentional modulations in favor of the attended stream, we still find that unattended speech is represented up to linguistic level in the regions that overlap with the reorienting

attention system. These linguistic representations of unattended speech might facilitate attentional reorienting and filtering during natural speech perception. Overall, our findings provide comprehensive insights on attentional mechanisms that underlie the ability to selectively listen to a desired speaker in noisy multi-speaker environments.

Funding

The work was supported in part by a National Eye Institute Grant (EY019684), by a European Molecular Biology Organization Installation Grant (IG 3028), by a TUBA GEBIP 2015 fellowship, and by a Science Academy BAGEP 2017 award.

Notes

The authors thank Jack L. Gallant, Wendy de Heer and Ümit Keleş for assistance in various aspects of this research. *Conflict of interest*: None.

References

Alho K, Medvedev SV, Pakhomov SV, Roudas MS, Tervaniemi M, Reinikainen K, Zeffirio T, Näätänen R. 1999. Selective tuning of the left and right auditory cortices during spatially directed attention. *Cogn Brain Res.* 7:335-341.

Alho K, Vorobyev VA, Medvedev SV, Pakhomov SV, Roudas MS, Tervaniemi M, Näätänen R. 2003. Hemispheric lateralization of cerebral blood-flow changes during selective listening to dichotically presented continuous speech. *Cogn Brain Res.* 17:201-211.

Alho K, Vorobyev VA, Medvedev SV, Pakhomov SV, Starchenko MG, Tervaniemi M, Näätänen, R. 2006. Selective attention to human voice enhances brain activity bilaterally in the superior temporal sulcus. *Brain Res.* 1075:142-150.

Alho K, Rinne T, Herron TJ, Woods DL. 2014. Stimulus-dependent activations and attention-related modulations in the auditory cortex: a meta-analysis of fMRI studies. *Hearing Res.* 307:29-41.

Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Statist Soc.* 57:289–300.

Bizley JK, Cohen YE. 2013. The what, where and how of auditory-object perception. *Nat Rev Neurosci.* 14:693-707.

Bregman AS. 1994. Auditory scene analysis: the perceptual organization of sound. MIT Press.

Broadbent D. 1958. Perception and communication. Pergamon Press

Brodbeck C, Presacco A, Simon JZ. 2018a. Neural source dynamics of brain responses to continuous stimuli: Speech processing from acoustics to comprehension. *NeuroImage.* 172:162-174

Brodbeck C, Hong LE, Simon JZ. 2018b. Rapid transformation from auditory to linguistic representations of continuous speech. *Curr Biology.* 28:3976-3983.

Boulenger V, Hoen M, Ferragne E, Pellegrino F, Meunier F. 2010. Real-time lexical competitions during speech-in-speech comprehension. *Speech Commun.* 52:246-253.

Bronkhorst AW. 2015. The cocktail-party problem revisited: early processing and selection of multi-talker speech. *Attent Percept Psychophys.* 77:1465-1487.

Carrasco M. 2011. Visual attention: the past 25 years. *Vision Res.* 51:1484-1525.

Cherry EC. 1953. Some experiments on the recognition of speech, with one and two ears. *J Acoust Soc Am.* 25:975–979.

Corbetta M, Shulman GL. 2002. Control of goal-directed and stimulus-driven attention in the brain. *Nat Rev Neurosci.* 3:201-215.

Corbetta M, Patel G, Shulman GL. 2008. The reorienting system of the human brain: from environment to theory of mind. *Neuron.* 58:306-324.

Da Costa S, van der Zwaag W, Miller LM, Clarke S, Saenz M. 2013. Tuning in to sound: frequency-selective attentional filter in human primary auditory cortex. *J Neurosci.* 33:1858-1863.

Dale AM, Fischl B, Sereno MI. 1999. Cortical surface-based analysis – I: Segmentation and surface reconstruction. *Neuroimage.* 9:179 –194.

Davis MH, Johnsrude IS. 2003. Hierarchical processing in spoken language comprehension. *J Neurosci.* 23:3423-3431.

de Heer WA, Huth AG, Griffiths TL, Gallant JL, Theunissen FE. 2017. The hierarchical cortical organization of human speech processing. *J Neurosci.* 37:6539–6557.

Degerman A, Rinne T, Salmi J, Salonen O, Alho K. 2006. Selective attention to sound location or pitch studied with fMRI. *Brain Res.* 1077:123–134.

Destrieux C, Fischl B, Dale A, Halgren E. 2010. Automatic parcellation of human cortical gyri and sulci using standard anatomical nomenclature. *Neuroimage.* 53:1–15.

DeWitt I, Rauschecker JP. 2012. Phoneme and word recognition in the auditory ventral stream. *Proc Natl Acad Sci.* 109:E505–E514.

Di Liberto GM, O’Sullivan JA, Lalor EC. 2015. Low-frequency cortical entrainment to speech reflects phoneme-level processing. *Curr Biol.* 25:2457–2465.

Ding N, Simon JZ. 2012a. Neural coding of continuous speech in auditory cortex during monaural and dichotic listening. *J Neurophysiol.* 107:78–89.

Ding N, Simon JZ. 2012b. Emergence of neural encoding of auditory objects while listening to competing speakers. *Proc Natl Acad Sci USA.* 109:11854–11859.

Ding N, Simon JZ. 2014. Cortical entrainment to continuous speech: functional roles and interpretations. *Front Hum Neurosci.* 8:311

Ding N, Pan X, Luo, C, Su N, Zhang W, Zhang J. 2018. Attention is required for knowledge-based sequential grouping: Insights from the integration of syllables into words. *J Neurosci.* 38:1178-1188.

Driver J. 2001. A selective review of selective attention research from the past century. *Brit J Psych.* 92:53-78

Elhilali M, Xiang J, Shamma SA, Simon JZ. 2009. Interaction between attention and bottom-up saliency mediates the representation of foreground and background in an auditory scene. *PLoS Biol* 7:e1000129.

Evans S, McGettigan C, Agnew ZK, Rosen S, Scott SK. 2016. Getting the cocktail party started: masking effects in speech perception. *J Cogn Neurosci.* 28:483-500.

Friederici AD. 2011. The brain basis of language processing: from structure to function. *Physiol Rev.* 91:1357-1392.

Fritz JB, Elhilali M, David SV, Shamma SA. 2007. Auditory attention—focusing the searchlight on sound. *Curr Opin Neurobiol.* 17:437-455.

Gao JS, Huth AG, Lescroart MD, Gallant JL. 2015. Pycortex: an interactive surface visualizer for fMRI. *Front Neuroinf.* 9.

Goutte C, Nielsen FA, Hansen K. 2000. Modeling the hemodynamic response in fMRI using smooth fir filters. *IEEE Trans Med Imag.* 19:1188–1201.

Greve DN, Fischl B. 2009. Accurate and robust brain image alignment using boundary-based registration. *Neuroimage.* 48:63–72.

Griffiths TD, Warren JD. 2004. What is an auditory object? *Nat Rev Neurosci.* 5:887-892.

Hickok G, Poeppel D. 2004. Dorsal and ventral streams: a framework for understanding aspects of the functional anatomy of language. *Cognition* 92:67-99

Hickok G, Poeppel D. 2007. The cortical organization of speech processing. *Nat Rev Neurosci.* 8:393–402.

Hill KT, Miller LM. 2010. Auditory attentional control and selection during cocktail party listening. *Cereb Cortex.* 20:583-590

Hink RF, Hillyard SA. 1976. Auditory evoked potentials during selective listening to dichotic speech messages. *Percept Psychophys.* 20:236–242.

Huth AG, De Heer WA, Griffiths TL, Theunissen FE, Gallant JL. 2016. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature.* 532:453-458.

Ikeda Y, Yahata N, Takahashi H, Koeda M, Asai K, Okubo Y, Suzuki H. 2010. Cerebral activation associated with speech sound discrimination during the diotic listening task: an fMRI study. *Neurosci Res.* 67:65-71.

Jäncke L, Buchanan TW, Lutz K, Shah NJ. 2001. Focused and nonfocused attention in verbal and emotional dichotic listening: an FMRI study. *Brain Lang.* 78:349–363.

Jäncke L, Specht K, Shah JN, Hugdahl K. 2003. Focused attention in a simple dichotic listening task: an fMRI experiment. *Cogn Brain Res*. 16:257-266.

Jenkinson M, Smith S. 2001. A global optimization method for robust affine registration of brain images. *Med Image Anal*. 5:143–156.

Johnson JA, Zatorre RJ. 2005. Attention to simultaneous unrelated auditory and visual events: behavioural and neural correlates. *Cereb Cortex*. 15:1609– 1620.

Kastner S, Pinsk MA. 2004. Visual attention as a multilevel selection process. *Cogn Affect Behav Neurosci*. 4:483–500.

Kerlin JR, Shahin AJ, Miller LM. 2010. Attentional gain control of ongoing cortical speech representations in a “Cocktail Party”. *J Neurosci*. 30:620–628.

Kidd G, Colburn HS. 2017. Informational masking in speech recognition. In: *The Auditory System at the Cocktail Party*. Springer. p. 75–109

Lavie N. 2005. Distracted and confused?: Selective attention under load. *Trends Cogn Sci*. 9:75-82.

Levelt WJ. 1993. *Speaking: from intention to articulation*. Cambridge (MA): MIT Press.

Lewis JL. 1970. Semantic processing of unattended messages using dichotic listening. *J Exp Psychol*. 85:225–228.

Li Y, Wang F, Chen Y, Cichocki A, Sejnowski T. 2018. The effects of audiovisual inputs on solving the cocktail party problem in the human brain: An fmri study. *Cereb Cortex*. 28:3623-3637.

Lipschutz B, Kolinsky R, Damhaut P, Wikler D, Goldman S. 2002. Attention-dependent changes of activation and connectivity in dichotic listening. *Neuroimage*. 17:643-656.

Maunsell JH, Treue S. 2006. Feature-based attention in visual cortex. *Trends Neurosci*. 29:317-322.

McDermott JH. 2009. The cocktail party problem. *Curr Biol*. 19:R1024-R102

Mesgarani N, Chang EF. 2012. Selective cortical representation of attended speaker in multi-talker speech perception. *Nature*. 485:233–U118.

Miller LM. 2016. *Neural Mechanisms of Attention to Speech*. In *Neurobiology of Language*. Academic Press. p. 503-514

Mitchell TM, Shinkareva SV, Carlson A, Chang KM, Malave VL, Mason RA, Just MA. 2008. Predicting human brain activity associated with the meanings of nouns. *Science*. 320:1191-1195.

Moray N. 1959. Attention in dichotic listening: Affective cues and the influence of instructions. *Q J Exp Psychol*. 11:56–60.

Nakai T, Kato C, Matsuo K. 2005. An fMRI study to investigate auditory attention: a model of the cocktail party phenomenon. *Magn Reson Med Sci*. 4:75-82.

Näätänen R, Kujala T, Winkler I. 2011. Auditory processing that leads to conscious perception: a unique window to central auditory processing opened by the mismatch negativity and related responses. *Psychophysio.* 48: 4-22.

Okada K, Rong F, Venezia J, Matchin W, Hsieh IH, Saberi K, Serences JT, Hickok, G. 2010. Hierarchical organization of human auditory cortex: evidence from acoustic invariance in the response to intelligible speech. *Cereb Cortex.* 20:2486–2495.

O'Craven KM, Downing PE, Kanwisher N. 1999. fMRI evidence for objects as the units of attentional selection *Nature.* 401:584-587.

O'Sullivan J, Herrero J, Smith E, Schevon C, Mckhann GM, Sheth SA, Mehta AH, Mesgarani N. 2019. Hierarchical Encoding of Attended Auditory Objects in Multi-talker Speech Perception. *Neuron.* 104:1195-1209.

Paltoglou AE, Sumner CJ, Hall DA. 2009. Examining the role of frequency specificity in the enhancement and suppression of human cortical activity by auditory selective attention. *Hear Res.* 257:106–118.

Petkov CI, Kang X, Alho K, Bertrand O, Yund EW, Woods DL. 2004. Attentional modulation of human auditory cortex. *Nat Neurosci.* 7:658–663.

Power AJ, Foxe JJ, Forde EJ, Reilly RB, Lalor EC. 2012. At what time is the cocktail party? A late locus of selective attention to natural speech. *Eur J Neurosci.* 35:1497–1503.

Pulvermüller F, Shtyrov Y. 2006. Language outside the focus of attention: the mismatch negativity as a tool for studying higher cognitive processes. *Prog Neurobio.* 79:49-71.

Puschmann S, Steinkamp S, Gillich I, Mirkovic B, Debener S, Thiel CM. 2017. The right temporoparietal junction supports speech tracking during selective listening: Evidence from concurrent EEG-fMRI. *J Neurosci.* 37:11505-11516.

Puschmann S, Baillet S, Zatorre RJ. 2019. Musicians at the cocktail party: neural substrates of musical training during selective listening in multispeaker situations. *Cereb Cortex.* 29:3253-3265.

Puvvada KC, Simon JZ. 2017. Cortical representations of speech in a multitalker auditory scene. *J Neurosci.* 37:9189-9196.

Rauschecker JP, Scott SK. 2009. Maps and streams in the auditory cortex: nonhuman primates illuminate human speech processing. *Nat Neurosci.* 12:718–724.

Rämä P, Relander-Syrjänen K, Carlson S, Salonen O, Kujala T. 2012. Attention and semantic processing during speech: an fMRI study. *Brain and Lang.* 122:114-119.

Regev M, Simony E, Lee K, Tan KM, Chen J, Hasson U. 2019. Propagation of information along the cortical hierarchy as a function of attention while reading and listening to stories. *Cereb Cortex.* 29:4017-4034.

Relander K, Rämä P, Kujala T. 2009. Word semantics is processed even without attentional effort. *J Cogn Neurosci.* 21:1511-1522.

Riecke L, Peters JC, Valente G, Kemper VG, Formisano E, Sorger B. 2017. Frequency-selective attention in auditory scenes recruits frequency representations throughout human superior temporal cortex. *Cereb Cortex*. 27:3002-3014.

Rinne T, Pekkola J, Degerman A, Autti T, Jääskeläinen IP, Sams M, & Alho K. 2005. Modulation of auditory cortex activation by sound presentation rate and attention. *Hum Brain mapping*. 26:94-99.

Rinne T, Balk MH, Koistinen S, Autti T, Alho K, Sams M. 2008. Auditory selective attention modulates activation of human inferior colliculus. *J Neurophysiol*. 100:3323-3327.

Rinne T. 2010. Activations of human auditory cortex during visual and auditory selective attention tasks with varying difficulty. *Open Neuroimage*. 4:187.

Rivenez M, Darwin CJ, Guillaume A. 2006. Processing unattended speech. *J Acoust Soc Am*. 119:4027-4040.

Sabri M, Binder JR, Desai R, Medler DA, Leitl MD, Liebenthal E. 2008. Attentional and linguistic interactions in speech perception. *Neuroimage*. 39:1444-1456.

Savitzky A, Golay MJ. 1964. Smoothing and differentiation of data by simplified least squares procedures. *Anal Chem*. 36:1627-1639.

Scott SK, Rosen S, Wickham L, Wise RJ. 2004. A positron emission tomography study of the neural basis of informational and energetic masking effects in speech perception. *J Acoust Soc Am*. 115:813–821.

Scott SK, Rosen S, Beaman CP, Davis JP, Wise RJS. 2009. The neural processing of masked speech: evidence for different mechanisms in the left and right temporal lobes. *J Acoust Soc Am*. 125:1737-1743.

Scott SK, McGettigan C. 2013. The neural processing of masked speech. *Hearing Res*. 303:58-66.

Seydell-Greenwald A, Greenberg AS, Rauschecker JP. 2014. Are you listening? Brain activation associated with sustained nonspatial auditory attention in the presence and absence of stimulation. *Hum Brain Mapp*. 35:2233–2252.

Shinn-Cunningham BG. 2008. Object-based auditory and visual attention. *Trends Cogn Sci*. 12:182-186.

Shinn-Cunningham BG, Best V. 2008. Selective attention in normal and impaired hearing. *Trends Amplif*. 12:283-299.

Shinn-Cunningham B, Best V, Lee AK. 2017. Auditory object formation and selection. In: *The Auditory System at the Cocktail Party*. Springer. p. 7–40.

Simon JZ. 2017. Human auditory neuroscience and the cocktail party problem. In: *The Auditory System at the Cocktail Party*. Springer. p. 169-197.

Smith SM. 2002. Fast robust automated brain extraction. *Hum. Brain Mapp*. 17:143–155.

Teder W, Kujala T, Näätänen R. 1993. Selection of speech messages in free-field listening. *Neuroreport*. 5:307-309

Treisman A. 1960. Contextual cues in selective listening. *Q J Exp Psychol*. 12:242-248

Treisman A. 1964. Monitoring and storage of irrelevant messages in selective attention. *J Verb Learn Verb Behav.* 3:449-459.

Vossel S, Geng JJ, Fink GR. 2014. Dorsal and ventral attention systems: distinct neural circuits but collaborative roles. *Neuroscientist.* 20:150-159

Wild CJ, Yusuf A, Wilson DE, Peelle JE, Davis MH, Johnsrude IS. 2012. Effortful listening: the processing of degraded speech depends critically on attention. *J Neurosci.* 32:14010-14021.

Wood N, Cowan N. 1995. The cocktail party phenomenon revisited: how frequent are attention shifts to one's name in an irrelevant auditory channel? *J Exp Psychol Learn Mem Cogn.* 21:255–260.

Woods DL, Stecker GC, Rinne T, Herron TJ, Cate AD, Yund EW, Liao I, Kang X. 2009. Functional maps of human auditory cortex: effects of acoustic features and attention. *PLoS ONE.* 4:e5183.

Woods DL, Herron TJ, Cate AD, Yund EW, Stecker GC, Rinne T, Kang X. 2010. Functional properties of human auditory cortical fields. *Front Syst Neurosci.* 4:155.

Xiang J, Simon J, Elhilali M. 2010. Competing streams at the cocktail party: exploring the mechanisms of attention and temporal integration. *J Neurosci.* 30:12084-12093.

Yuan J, Liberman M. 2008. Speaker identification on the SCOTUS corpus. *J Acoust Soc Am.* 123:3878.

Zion Golumbic EM, Ding N, Bickel S, Lakatos P, Schevon CA, Mckhann GM, Goodman RR, Emerson R, Mehta AD, Simon JZ. 2013. Mechanisms underlying selective neuronal tracking of attended speech at a 'cocktail party'. *Neuron.* 77:980-991.

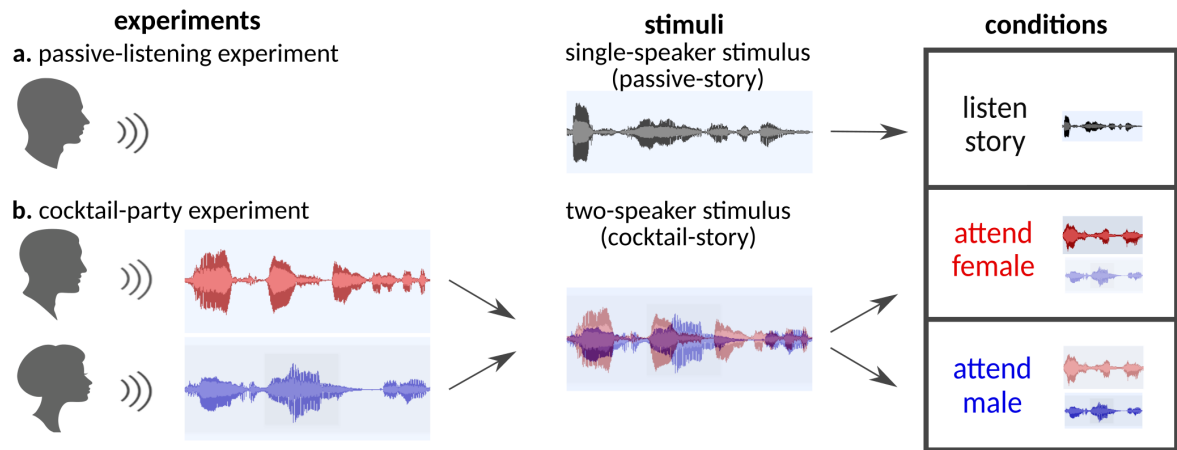


Figure 1: Experimental design. **a. Passive-listening experiment.** 10 stories from Moth-Radio-Hour were used to compile a single-speaker stimulus set (passive-stories). Subjects were instructed to listen to the stimulus vigilantly without any explicit task in the passive-listening experiment. **b. Cocktail-party experiment.** A pair of stories told by individuals of different genders were selected from the passive-story stimulus set and overlaid temporally to generate a two-speaker stimulus set (cocktail-stories). Subjects were instructed to attend either to the male or female speaker in the cocktail-party experiment. The same cocktail-story was presented twice in separate runs while the target speaker was varied. Attention condition was fixed within runs and it alternated across runs.

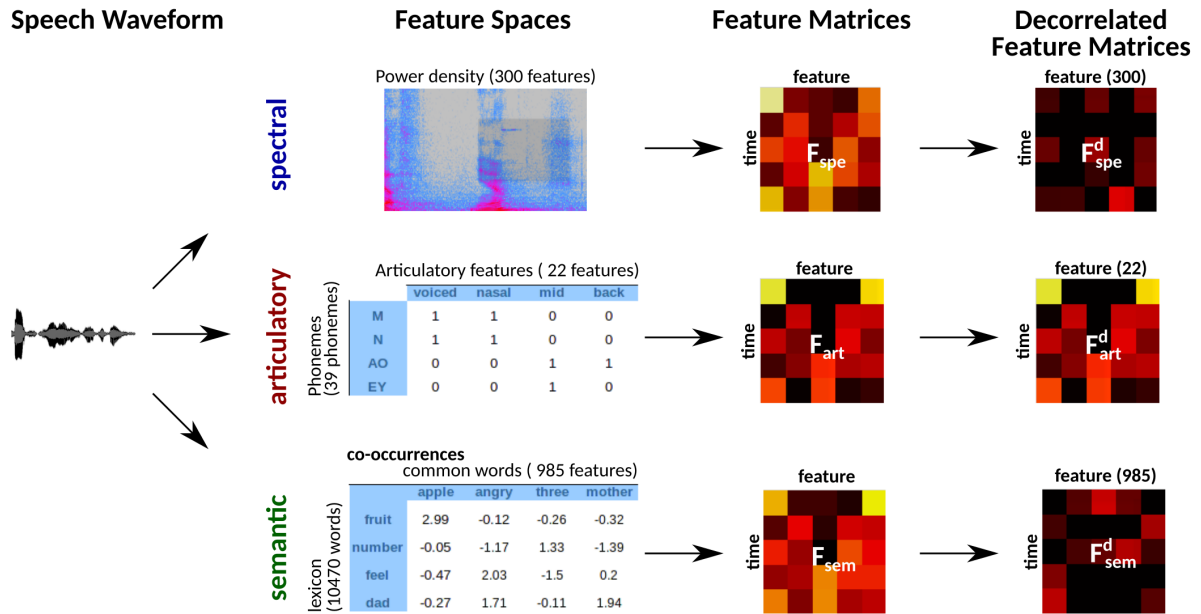


Figure 2: Multi-level speech features. Three distinct feature spaces were constructed to represent natural speech at multiple levels: namely spectral, articulatory, and semantic spaces. Speech waveforms were projected separately on these spaces to form stimulus matrices for each feature space. As such, the spectral feature matrix captured the spectral power density of the stimulus in 300 50-Hz bands between 0 and ~15 kHz. The articulatory feature matrix captured the mapping of each phoneme in the stimulus to 22 binary articulation features. The semantic feature matrix captured the statistical co-occurrences of each word in the stimulus with 985 common words in English. Each feature matrix was Lanczos-filtered at a cutoff frequency of 0.25 Hz and downsampled to 0.5 Hz to match the sampling rate of fMRI. Natural speech might contain intrinsic stimulus correlations among spectral, articulatory, and semantic features. If substantial, these correlations can in turn bias estimates of feature selectivity. To prevent potential biases, we decorrelated the three feature matrices examined here via Gram-Schmidt orthogonalization (see Methods). Taking the articulatory feature matrix as a reference, articulatory features were regressed out of the spectral feature matrix, and both articulatory and spectral features were regressed out of the semantic feature matrix. While spectral features are presumably lower level than articulatory features, the articulatory feature matrix with orders of magnitude lower number of features was taken as reference to minimize accumulation of numerical errors.

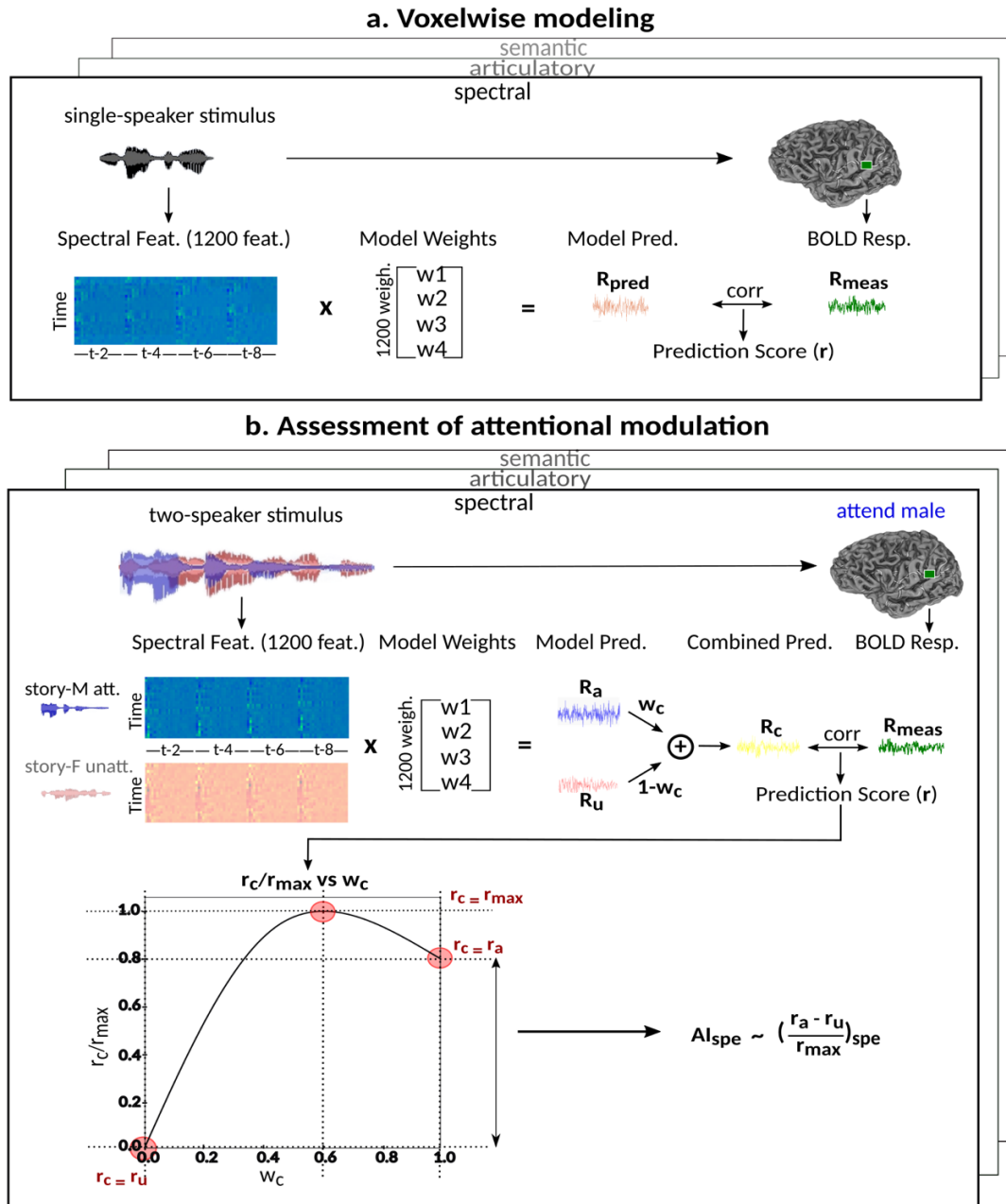


Figure 3: Analysis procedures. **a. Voxelwise modeling.** Voxelwise models were fit in individual subjects using passive-data. To account for the temporal characteristics of the hemodynamic response, a linearized four-tap finite impulse response (FIR) filter was used; and each model feature was represented as four features in the stimulus matrix to account for their delayed effects in BOLD responses at 2-8 sec. Model weights were found using L2-regularized linear regression. BOLD responses were predicted based on fit voxelwise models on held-out passive-data. Prediction scores were taken as the Pearson's correlation between predicted and measured BOLD responses. For a given subject, speech-selective voxels were taken as the union of voxels significantly predicted by spectral, articulatory, or semantic models ($q(\text{FDR}) < 10^{-4}$, t-test). **b. Assessment of attentional modulation.** Passive-models for single voxels were tested on cocktail-data to quantify attentional modulations in selectivity for corresponding model features. In a given run, one of the two speakers in a cocktail-story was attended while the other speaker was ignored. Separate response predictions were obtained using the isolated story stimuli for the attended speaker and for the unattended speaker. Since a voxel can represent information from both attended and unattended stimuli, a linear combination of these predicted responses was considered with varying combination weights (w_c in $[0, 1]$). BOLD responses were predicted based on each combination weight separately. Three separate prediction scores were calculated based on only the attended stimulus ($w_c=1$), based on only the unattended stimulus ($w_c=0$), and based on the optimal combination of the two stimuli. A model-specific attention index, (AI_m) was then computed as the ratio of the difference in prediction scores for attended versus unattended stories to the prediction score for their optimal combination (see Methods).

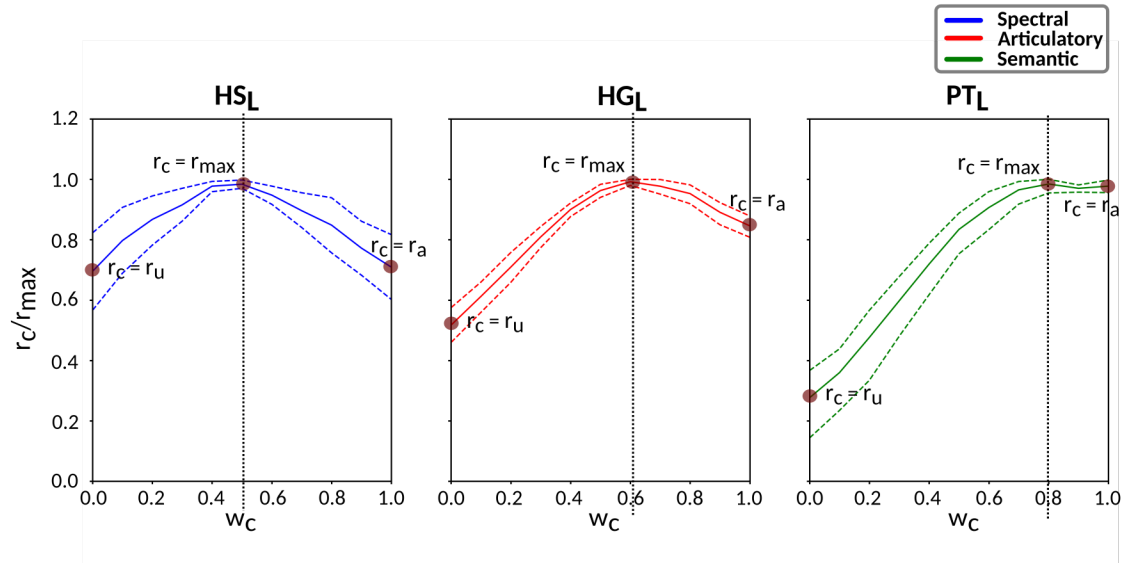


Figure 4: Predicting responses measured during the cocktail-party task. Passive-models were tested during the cocktail-party task by predicting BOLD responses in the cocktail-data. Since a voxel might represent information from both attended and unattended stimuli, response predictions were expressed as a convex combination of individual predictions for the attended and unattended story within each cocktail-story. Prediction scores were computed as the combination weights (w_c) were varied in [0 1] (see Methods). Prediction scores for a given model were averaged across speech-selective voxels within each ROI (r_c). The normalized scores of spectral, articulatory, and semantic models are displayed as a function of the combination weight in several representative ROIs (HS, HG and PT, respectively). Solid and dashed lines indicate mean and %95 confidence intervals across subjects. Scores based on only the attended story (r_a), based on only the unattended story (r_u), and based on the optimal combination (r_{max}) are marked with circles. For the *spectral model* in left HS, the optimal combination equally weighs attended and unattended stories, while no significant difference exists between r_a and r_u ($p > 0.05$, bootstrap test). This result implies that attention does not affect spectral representations in left HS. For the *articulatory model* in left HG, r_a is larger than r_u ($p < 10^{-4}$), while r_{max} is greater than r_a ($p < 10^{-2}$). This finding suggests that attention moderately impacts articulatory representations in left HG such that articulatory representations of the unattended story are still maintained. For the *semantic model* in left PT, r_a is greater than r_u ($p < 10^{-4}$), while no significant difference exists between r_{max} and r_a ($p > 0.05$). This finding indicates that attention strongly biases semantic representations in left PT towards the attended stimulus. These representative results imply that attention may have divergent effects at various levels of speech representations across cortex.

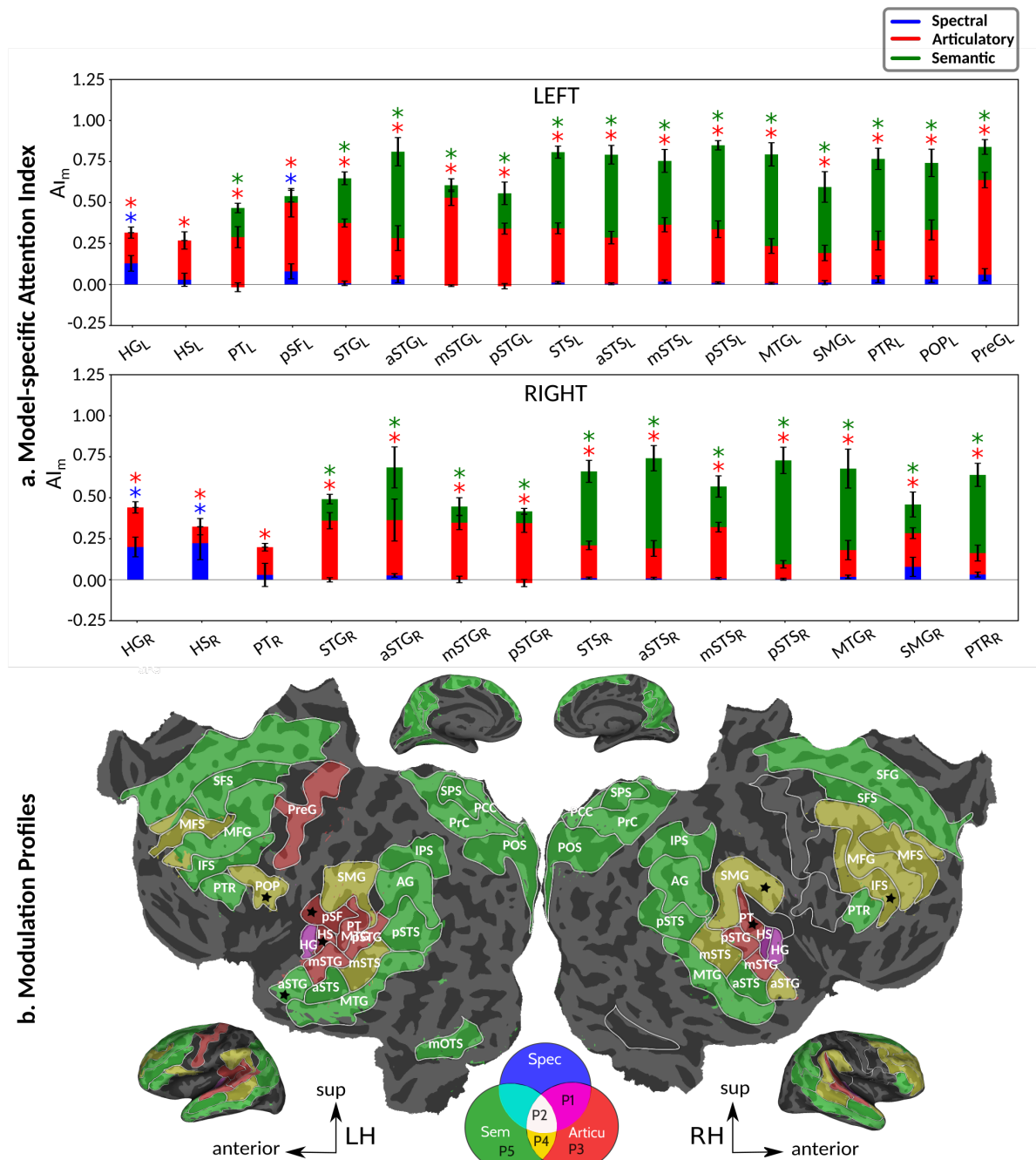


Figure 5: Attentional modulations for multi-level speech features. **a. Model-specific attention indices.** To quantify modulatory effects on selectivity for multi-level speech features, model-specific attention index (AI_m) was computed based on the difference in model prediction scores when the stories were attended versus unattended (see Methods). AI_m is in the range of $[-1,1]$, where a positive index indicates attentional modulation in favor of the attended stimuli and a negative index indicates modulation in favor of the unattended stimuli. A value of zero indicates no modulation. Blue, red, and green bar plots display spectral, articulatory, and semantic attention indices, respectively (mean \pm sem across subjects). The sum of the three indices gives the overall modulation; hence, the relative values of indices represent the composition of overall modulation in terms of the underlying feature spaces (see Methods). Significant modulations are marked with *, colored according to the corresponding metric ($p < 0.05$, bootstrap test; see legend; also see Supplementary Table S2b for detailed results). Only ROIs in the perisylvian cortex where a speech-related fronto-temporo-parietal network resides (Friederici 2011) are displayed (see Supplementary Fig. S6 for non-perisylvian ROIs). ROIs in the LH and RH are shown in top and bottom panels, respectively. POP_R and $PreG_R$ that did not have consistent speech selectivity in individual subjects were excluded (see Methods) These results show that selectivity modulations distribute broadly across cortex at the linguistic level (articulatory and semantic), yet modulations at the acoustic level (spectral) are primarily constrained to early auditory cortex. **b. Attentional modulation profiles.** Modulation profiles of cortical ROIs averaged across subjects are displayed on the flattened cortical surface of a representative subject (S4). Medial and lateral views of the inflated hemispheres are also shown above and below the flatmap. White lines encircle ROIs that are found based on an automatic atlas-based cortical parcellation.

Labels of ROIs are shown (see Methods for ROI abbreviations). This analysis only included ROIs with consistent selectivity for speech features in each individual subject (see Methods). Colors indicate the modulation profiles based on dominant AI_m (see legend). The distribution of model-specific attentional indices (AI_m) across spectral, articulatory, and semantic models were examined and compared to the selectivity profiles of the ROIs (see Supplementary Fig. S2 for intrinsic selectivity profiles). Among ROIs with intrinsic selectivity profile P1, AI_{art} and AI_{spe} are dominant with exception in left HS and pSF, and right PT where AI_{art} is dominant. For P2 and P4, AI_{art} and AI_{sem} are dominant with exception in aSTG where only AI_{sem} is dominant. For P3, AI_{art} is dominant. For P5, AI_{sem} is dominant with exception in right IFS and left POP where both AI_{sem} and AI_{art} are dominant (see Supplementary Table S2b for detailed results). Overall, the intrinsic selectivity and attentional modulation profiles display largely matching distributions across cortex (exceptions are marked with *; see Supplementary Fig. S7 for statistical assessments). Exceptions are mainly attributed to the lack of spectral modulation outside early auditory cortex. These results support the view that attention enhances speech representations to favor target features that each brain area is intrinsically selective for, albeit this effect is relatively weaker at the spectral level.

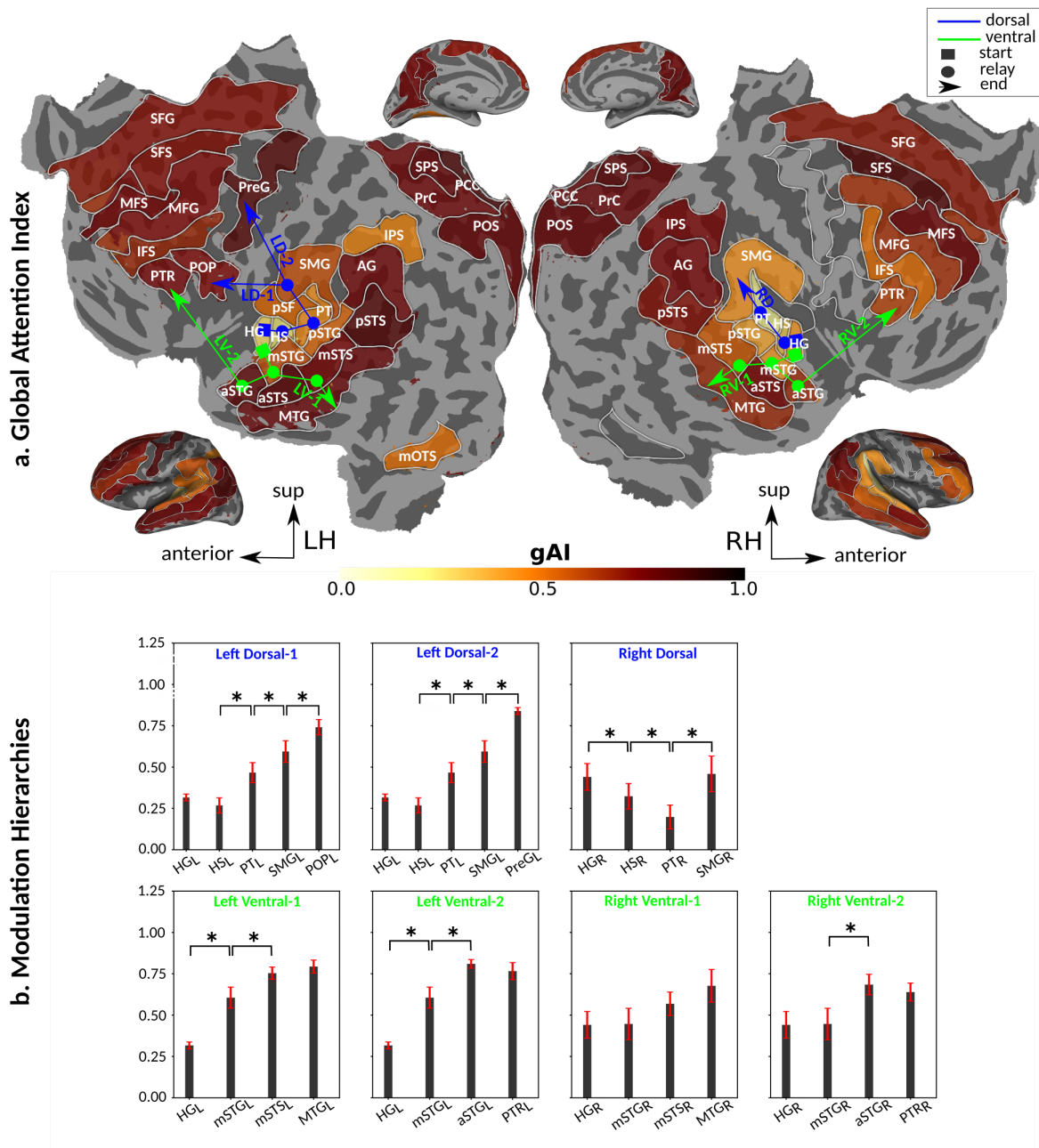


Figure 6: Global attentional modulation. a. Global attention index and auditory pathways. To quantify overall modulatory effects on selectivity across all examined feature levels, global attentional modulation (gAI) was computed by summing spectral, articulatory, and semantic attention indices (see Methods). gAI is in the range of $[-1,1]$, where a positive index indicates attentional modulation in favor of the attended stimuli and a negative index indicates modulation in favor of the unattended stimuli. A value of zero indicates no modulation. Colors indicate gAI averaged across subjects (see legend; see Supplementary Fig. S8 for bar plots of gAI across cortex) To illustrate the gradients in gAI across the hierarchy of speech processing, dorsal and ventral pathways are shown with blue and green lines, respectively. Squares mark the region where the pathway begins; arrows mark the region where the pathway ends; and circles mark the relay regions in between. Dorsal pathway comprises three trajectories: left dorsal-1 (LD-1: $HG_L \rightarrow HS_L \rightarrow PT_L \rightarrow SMG_L \rightarrow POP_L$), left dorsal-2 (LD-2: $HG_L \rightarrow HS_L \rightarrow PT_L \rightarrow SMG_L \rightarrow PreG_L$) and right dorsal (RD: $HG_R \rightarrow HS_R \rightarrow PT_R \rightarrow SMG_R$). Ventral pathway comprises four trajectories: left ventral-1 (LV-1: $HG_L \rightarrow mSTG_L \rightarrow mSTS_L \rightarrow MTG_L$), left ventral-2 (LV-2: $HG_L \rightarrow mSTG_L \rightarrow aSTG_L \rightarrow PTR_L$), right ventral-1 (RV-1: $HG_R \rightarrow mSTG_R \rightarrow mSTS_R \rightarrow MTG_R$) and right ventral-2 (RV-2: $HG_R \rightarrow mSTG_R \rightarrow aSTG_R \rightarrow PTR_R$). **b. Modulation hierarchies.** Gradients in gAI across left dorsal-1 and dorsal-2, right dorsal, left ventral-1 and ventral-2, and right ventral-1 and ventral-2 trajectories are shown in subfigures labeled accordingly. Bar plots display gAI in each ROI (mean \pm sem across subjects). Only ROIs within a given trajectory are included in the corresponding subfigure. Significant differences in gAI between consecutive ROIs along the trajectory are marked with brackets ($p < 0.05$, bootstrap test; see Table S2d for detailed results). Significant gradients in gAI are: $gAI_{HS} < gAI_{PT} < gAI_{SMG} < gAI_{POP}$ in left dorsal-1, $gAI_{HS} < gAI_{PT} < gAI_{SMG} < gAI_{PreG}$ in left dorsal-2, $gAI_{PT} < gAI_{SMG}$ and $gAI_{HG} > gAI_{HS} > gAI_{PT}$ in right dorsal, $gAI_{HG} < gAI_{mSTG} < gAI_{mSTS}$ in left ventral-

1, $gAI_{HG} < gAI_{mSTG} < gAI_{aSTG}$ in left ventral-2, and $gAI_{mSTG} < gAI_{aSTG}$ in right ventral-2. These results indicate that in the left hemisphere, gAI gradually increases from early auditory regions to higher-order regions across the dorsal and ventral pathways. Furthermore, the gradients in gAI mostly overlap with the gradients in CI ($r > 0.84$, $p < 10^{-4}$, bootstrap test; see Supplementary Fig. S5 for gradients in CI). Similar patterns are also observed in the right hemisphere, although the gradients in gAI are less consistent. These results suggest that, with growing representational complexity, attentional modulation also grows stronger across the cortical hierarchy of speech.

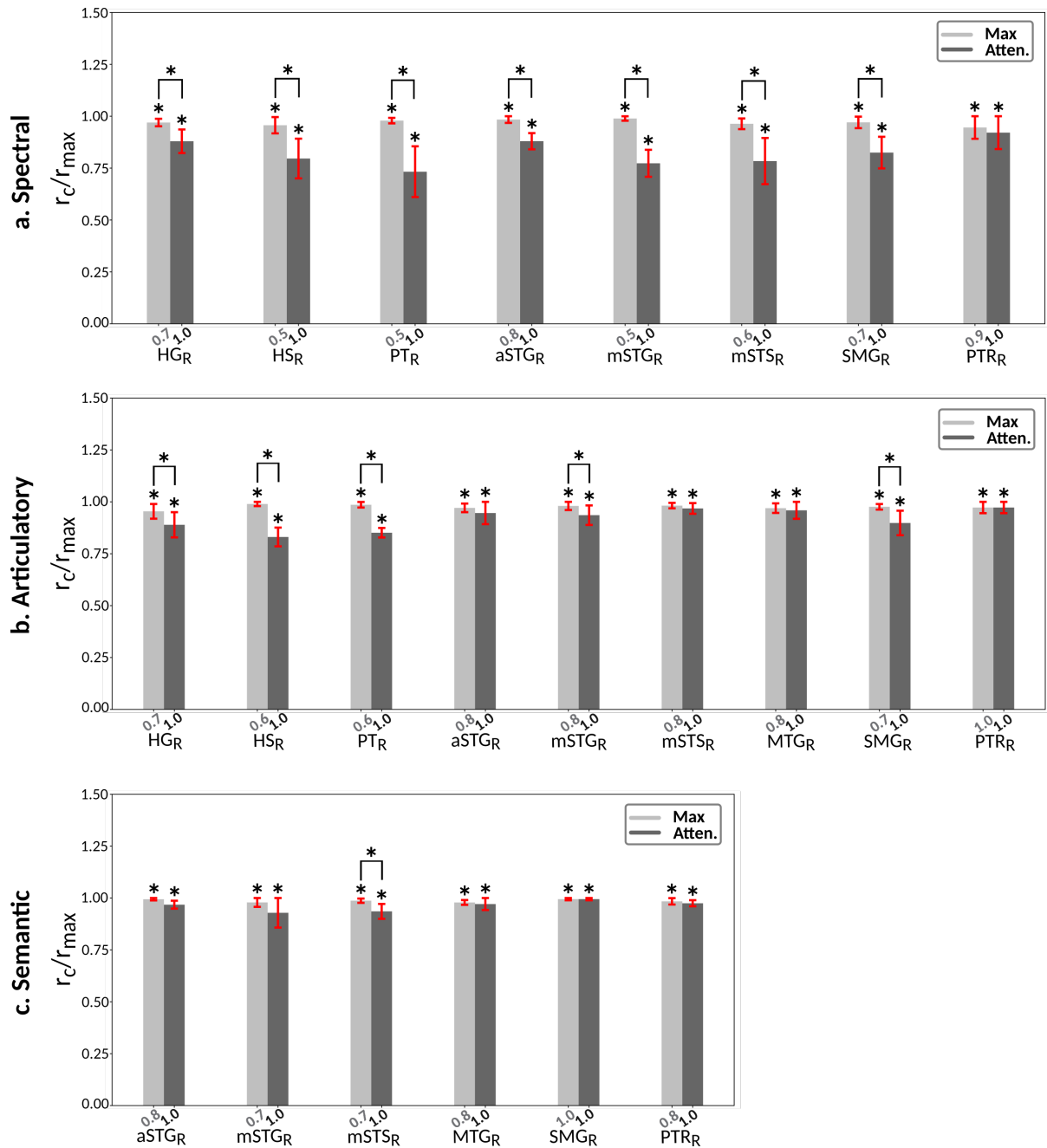


Figure 7: Representation of unattended speech. Passive models were tested on cocktail-data to assess representation of unattended speech during the cocktail-party task. Prediction score of a combination model comprising the features of both attended and unattended stories (r_{max} : optimal convex combination) and prediction score of an individual model comprising only the features of the attended story (r_a) were computed (see Methods). Significant differences between the two prediction scores indicate that BOLD responses within the ROI carry significant information on unattended speech. Bar plots display normalized prediction scores (mean \pm sem across subjects) for the optimal convex combination (light gray) and the attended story (gray). Significant scores are marked with * ($p < 10^{-4}$, bootstrap test). Brackets indicate significant differences between the two scores ($p < 0.05$). Prediction scores are only displayed for ROIs in the right dorsal and ventral stream, with significant selectivity for corresponding model features (see Fig. S9 for left hemisphere). **a. Spectral representation.** Spectral representations of unattended speech extend up to SMG along the dorsal stream (HG \rightarrow HS \rightarrow PT \rightarrow SMG) and up to mSTS and aSTG along the ventral stream (HG \rightarrow mSTG \rightarrow mSTS \rightarrow MTG and HG \rightarrow mSTG \rightarrow aSTG \rightarrow PTR). **b. Articulatory representation.** Articulatory representations of unattended speech extend up to SMG along the dorsal stream, and up to mSTG along the ventral stream. **c. Semantic representation.** Semantic features are represented only in mSTS. These results suggest that processing of unattended speech is not constrained at spectral level but extends to articulatory and semantic level.