

# Interaural time and level differences contribute differently to source segregation and spatial selection

Benjamin N. Richardson, Sahil Luthra, Jana M. Kainerstorfer, Barbara G. Shinn-Cunningham, Christopher A. Brown

## Abstract

Interaural time and level differences are crucial in sound localization, yet their contributions to sound source segregation and spatial selection remain underspecified. Here, participants completed a spatial auditory selective attention task while we measured hemodynamic activity in the prefrontal cortex and superior temporal gyrus using functional near-infrared spectroscopy. Participants listened to a target sound stream and a simultaneous spatially separated speech stream or white noise masker. Sound streams were spatialized with either 50  $\mu$ s ITDs, 500  $\mu$ s ITDs, naturally occurring ILDs from a non-individualized HRTF, or broadband 10 dB ILDs. Behavioral results revealed a stronger effect of spatial cues when the masker was speech. Error patterns differed in the two difficult conditions, small ITDs and natural ILDs: small ITDs produced lower hit rates, while naturally occurring ILDs produced higher false alarm rates. Small ITDs led to greater activity in prefrontal cortex and activity in superior temporal gyrus that was lateralized, greater in the hemisphere contralateral to attentional focus, consistent with previous reports. These results suggest that natural ILDs alone support source segregation even if they are insufficient to cause large shifts in perceived lateralization, explaining high false alarm rates (confusions between target and distractor words). In contrast, small ITDs alone may be insufficient to segregate competing sources, leading to low hit and false alarm rates. Together, these results reveal differences in how ITDs and ILDs contribute to auditory scene analysis and spatial attention.

## Introduction

In many typical sound environments, the signals reaching a listener's ears comprise a rich mixture of concurrent, spectro-temporally overlapping but independent sounds, often arriving from different directions (Cherry, 1953). To understand a sound in such settings, the brain must achieve two interrelated tasks: 1) *segregate* the sources in the mixture, forming distinct perceptual objects by determining what acoustic energy belongs to what source and 2) *select*

which object to process in detail by focusing selective attention to filter out competing sources (Noyce et al., 2023; Shinn-Cunningham, 2008a). Many acoustic factors influence both segregation and selection, including pitch, temporal envelope cues, and spatial features, among others (Middlebrooks and Waters, 2020; Moore and Gockel, 2002). The current work examines how spatial cues influence segregation and spatial selection.

Spatial cues are relatively weak compared to spectrotemporal cues in grouping together simultaneous sound elements, but they play a larger role in “streaming” sounds across silent gaps (Ihfeldt and Shinn-Cunningham, 2008; Middlebrooks, 2017; Moore and Gockel, 2002). Perceived spatial separation can also guide spatial attention to a target in a mixture (Arbogast et al., 2002; Freyman et al., 1999; Kidd et al., 2010). Yet spatial cues have little influence if competing sounds have different and distinct spectrotemporal structure, such as when a woman speaks over a background of radio static. When spectrotemporal features alone are enough to support segregation and selection, spatial separation between competing sources has little effect on psychophysical outcomes (Culling and Summerfield, 1995; Darwin, 2007; Kubovy, 1988; Kubovy and Van Valkenburg, 2001; Noyce et al., 2023). Thus, the impact of spatial separation on auditory attention depends critically on the spectrotemporal properties of competing sounds.

The influence of spatial cues on segregation and selection helps explain spatial release from masking (SRM), wherein listeners are better at detecting and understanding a target in the presence of a spatially separated masker compared to when they are co-located (Freyman et al., 1999; Middlebrooks and Waters, 2020). The most important and robust spatial cues for SRM are interaural time and level differences (ITDs and ILDs). These cues help listeners to selectively attend to a target, though their specific contributions to segregation versus selection remain underspecified.

ITDs and ILDs have distinct effects on the signals reaching the ears, which may impact how they contribute to SRM. Acoustically, ITDs do not alter the target-to-masker ratio (TMR) in the mixture of signals reaching the ears; thus, any benefits from ITDs must arise from explicit spatial computations within the auditory pathway. In contrast, ILDs affect the TMR in the mixture of signals reaching each ear; typically, when a target and a masker have different ILDs, the TMR will be larger at the ear nearer the target, producing a “better ear” effect (Glyde et al, 2013). Perceptually, this TMR boost can improve target understanding even in the absence of perceived spatial separation or in monaural listening conditions by simply reducing energetic masking. Physiologically, a source containing ILDs is preferentially represented in the hemisphere contralateral to the ear receiving the more intense signal (Higgins et al., 2017). A direct consequence of this cortical lateralization is that the neural representations of competing sounds with opposing ILDs show inherently less overlap than representations of sources containing only ITDs. This lateralization may support perceptual source segregation even in the absence of explicit spatial effects. These differences suggest that ITDs and ILDs may support SRM through different mechanisms.

The extent to which ILDs support SRM will also depend on the spectrotemporal properties of competing sounds. Natural ILDs are substantial only in high frequencies (Yost and Dye, 1988). Because high-frequency components of speech group with their corresponding low-frequency components (Best et al., 2007; Bregman, 1990; Darwin, 1997), natural ILDs may nevertheless help with perceptual segregation of competing speech. Yet when paired with uninformative or zero ITDs, the perceived location of speech containing only natural ILDs and no usable ITDs will be nearer midline, since ITDs at low and mid frequencies dominate lateralization (Ellinger et al., 2017; Glyde et al., 2013; Wightman and Kistler, 1992). As a result, competing sources with only natural, opposing ILDs may be less likely to be perceived as coming from distinct locations, making source selection difficult. In contrast, broadband ILDs—applied uniformly across frequencies—can provide strong grouping and separation cues.

Magnifying ILDs to larger-than-natural magnitudes can enhance SRM in both normal-hearing listeners and cochlear implant users (Brown, 2014; Richardson et al., 2025). Without ILD magnification, CI users typically do not show significant SRM, likely because cochlear implants do not preserve ITDs and compress ILDs (Brown, 2018; Dorman et al., 2014; Gray et al., 2021). The current study seeks to clarify how ILDs contribute to SRM in normal hearing listeners in order to begin to identify how enhanced ILDs benefit CI users.

Functional neuroimaging studies reveal that spatial attention engages a visuo-spatial control network that includes prefrontal cortex (PFC; (Kong et al., 2014; Michalka et al., 2015; Noyce et al., 2022) and enhances contralateral bias in the representation of sound in auditory cortex (Alho et al., 1998, 1999, 2003; Ciaramitaro et al., 2007; Jäncke and Shah, 2002; Yang and Mayer, 2014). Furthermore, previous functional near-infrared spectroscopy (fNIRS) studies show that PFC activity varies with cognitive demands, reflecting individual differences in speech intelligibility (Lawrence et al., 2018; Zhou et al., 2022a, 2022b), perceived task difficulty (Zhou et al., 2022a, 2022b), and spatial separation between target and masker sounds (Zhang et al., 2021a). Measures of effort, including hemodynamic activity in PFC, show a non-monotonic relationship with task difficulty: reduced responses can reflect either very low or very high difficulty (Lawrence et al., 2018). In a spatial selective attention task similar to that used in the present study, PFC showed less activity when target and masker were dichotic than when lateralized with relatively small ITDs (Zhang et al., 2021a). This result could reflect reduced demands on spatial attention when target and masker locations are more perceptually distinct (infinite ILDs compared to small ITDs); however, it could also reflect reduced effort expended to segregate sources lateralized with ILDs than with ITDs. Specifically, because neural representations are automatically lateralized for sources with ILDs, computational demands of segregation and selection may be reduced, resulting in decreased activation in PFC.

Here, we used fNIRS to estimate neural activity in both PFC and auditory cortex while participants performed a spatial selective attention task. We compared behavioral and neural responses with small ITDs, large ITDs, natural ILDs, and broadband ILDs under conditions where the masker was either white noise or competing speech. We hypothesized that when the masker was speech rather than noise, spatial cues would strongly influence patterns of behavior because of the relatively higher spectrotemporal similarity. We expected greater spatial-cue

benefits with speech maskers than with noise maskers, and that these benefits would be reflected in both improved performance and reduced PFC activity. We also expected that, relative to small ITDs and natural ILDs, large ITDs and broadband ILDs would lead to improved behavioral performance and reduced PFC engagement, reflecting less difficult listening conditions. We hypothesized that lateralization of speech with natural ILDs would not strongly support spatial separation of target and masker, leading to poorer performance than with broadband ILDs. Finally, we predicted that STG activity would lateralize contralateral to the attended direction, consistent with previous reports. The current investigation provides insights into how ITDs and ILDs contribute to SRM and highlights how the different spatial cues shape both perceptual performance and cortical engagement.

## Materials & Methods

### Participants

Thirty native English speaking listeners participated in this experiment (27 female, 3 male, age  $22.6 \pm 4.45$  SD). All had normal hearing, defined as having pure-tone audiometric thresholds of 20 dB HL or better at octave frequencies from 250 to 8 kHz. All listeners gave written informed consent prior to participating in the study. All testing was administered according to the guidelines of the Institutional Review Board of the University of Pittsburgh.

### Overview of the Task

Subjects completed a color-word detection task in a selective attention paradigm with two competing, spatially separated sound streams. Figure 1 shows a schematic of target and masker streams. The target stream consisted of temporally jittered words spoken by a single male talker. The competing stream was either speech from the same talker or white noise. Within an experimental block, target and masker sequences were spatialized such that they were perceived to originate from opposing hemifields, one on the left and the other on the right, assigned randomly on each block.

For the speech masker condition, target and masker sequences each consisted of 18 word tokens, which were presented in pairs (one target token and one masker token). A word token pair occurred every 600 ms. Within each token pair, the two 300-ms-long tokens partially overlapped in time, with the second one starting 150 ms after the first. To ensure that listeners relied on spatial attention rather than temporal regularity, which token in a pair occurred first was randomly selected, independently, for each pair in a block. A 150-ms-long silent gap followed the second token before the next pair began. For the noise masker condition, masker word tokens were replaced by white noise bursts.

The diagram illustrates the experimental paradigm over a 12.8-second period. It shows the timing of the target word, speech mask, and noise mask. The target words are 'bag', 'box', 'blue', 'pen', and 'toy'. The speech mask words are 'shoe', 'card', 'box', and 'white'. The noise mask is a 300 ms duration. The sequence starts at 0 s (Cue) and ends at 12.8 s. The 'blue' target and 'card' mask pair are highlighted with 'Lead syllable' and 'Lag syllable' labels. The sequence ends with a 'Rest' period.

**Figure 1.** Schematic of a single block. Listeners first heard a cue word “bag” from the target stream. Then, they responded by button press to color words in the target stream (top row). The target stream was presented against either a speech or noise masker at the same intensity (i.e., with a 0-dB target-to-masker ratio). Each block was followed by 13-18 seconds rest.

## Stimuli

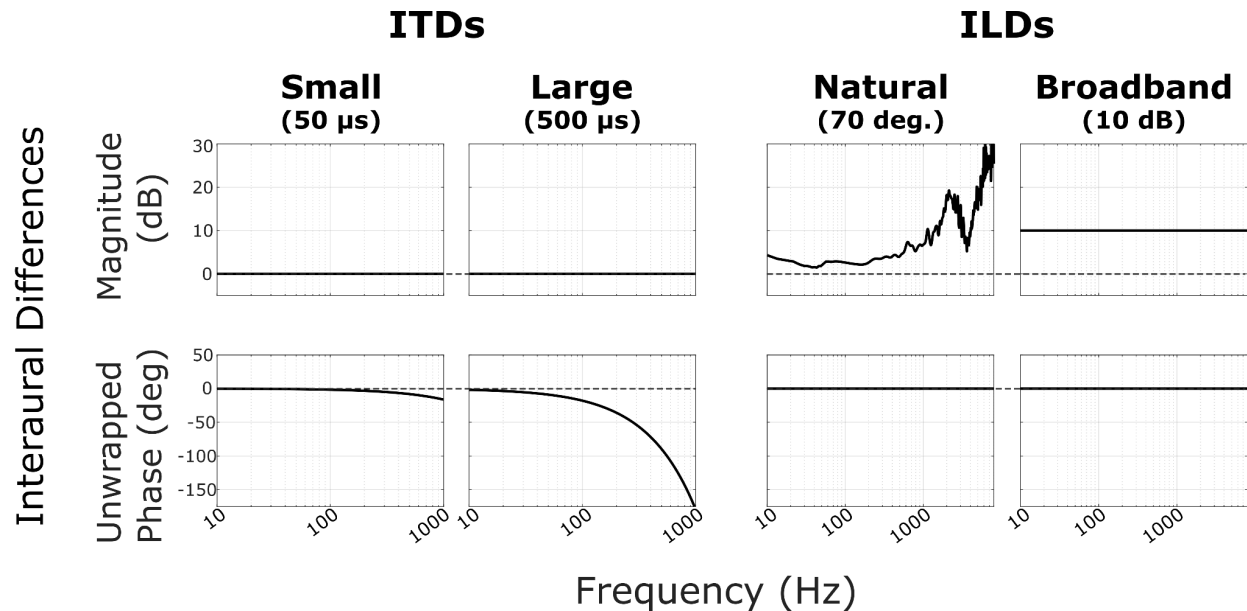
All stimulus processing was completed in the Python programming language (Van Rossum and Drake, 2009), with signal processing functions that were either from the *scipy* package (Virtanen et al., 2020) or custom written.

Target tokens were American English words. These speech stimuli were constructed from a list of 16 words, recorded in isolation by a single male talker (Kidd et al., 2008). The words were selected from a set of 12 object words: *<hat, bag, card, chair, desk, glove, pen, shoe, sock, spoon, table, toy>* and a set of four color words: *<red, white, blue, green>*. All word tokens were edited to have a duration of 300 ms by using a waveform editor to isolate each token and the soundstretch synchronous-overlap-and-add algorithm to adjust duration (“SoundStretch Audio Processing Utility,” n.d.). Finally, tokens were equated in root-mean-square amplitude.

Target and masker word sequences were generated by randomly choosing 18 words, with replacement, from the object and color word sets. Each target stream included a minimum of 3 and a maximum of 5 color words; the other words were all from the object set. The temporal position of color words was restricted such that color words could not occur within 2 indices of each other in either stream.

Two types of maskers were used: speech (second row, Figure 1) and noise (third row, Figure 1). Speech masker sequences were generated identically to target sequences, but with the additional constraint that within each target-masker pair, the masker word differed from the target word. Note that color words could occur in the speech masker sequence, but the listener was asked to ignore these “foil” color words. Noise masker streams comprised independent 300-ms white noise bursts. Noise bursts were subjected to the same timing constraints as speech masker word tokens (Fig. 1). Color words across both target and speech masker streams were constrained to ensure that color words could not occur in successive word token pairs. The presentation level of each stream was set to 65 dB SPL resulting in a 0 dB target-to-masker ratio prior to applying spatial cues.

The four different spatialization conditions are illustrated in Figure 2. To each stream (target and masker, spatialized to opposite lateral locations), we either applied  $\pm 50 \mu\text{s}$  ITD (“Small ITD”),  $\pm 500 \mu\text{s}$  ITD (“Large ITD”), a naturally occurring ILD corresponding to 70 degrees azimuth (“Natural ILD”), or a 10 dB broadband ILD (“Broadband ILD”). To apply ITDs, the signal in the ear contralateral to the apparent source location was delayed in time by the appropriate number of samples (corresponding to 50  $\mu\text{s}$  or 500  $\mu\text{s}$ ). In the natural ILD condition, the magnitude spectrum from an HRTF for 70° azimuth, measured on a KEMAR manikin, was applied to spatialize sources to the right (Gardner and Martin, 1995). For sources to the left, the left and right magnitude spectra from this HRTF were swapped. The 70° location was chosen because it yielded the greatest average broadband ILD magnitude, equal to about 16 dB. In the broadband ILD condition, the signal in the contralateral ear was attenuated by 10 dB (equally across all frequencies). After spatialization, target and masker streams, one from the left and one from the right, were summed to create left and right channels.



**Figure 2.** Spatialization conditions. The left (L) and right (R) ear signals for a given stream were spatialized with a 50  $\mu$ s ITD (“Small ITD”), 500  $\mu$ s ITD (“Large ITD”), naturally occurring ILD (HRTF at 70° azimuth, “Natural ILD”), or a 10 dB attenuation (“Broadband ILD”). Each subplot shows the magnitude (top row) and unwrapped phase (bottom row) properties of the interaural differences used. ITDs were implemented by pure time delay of the contralateral ear. Note that for the ITD plots, the unwrapped phase plots show a characteristic falloff with frequency proportional to  $1/f$ , which correspond to group delays of 50  $\mu$ s and 500  $\mu$ s for small and large ITDs, respectively.

## Experimental Procedure

Participants were seated in an anechoic chamber, fitted with Etymotic ER-3a insert phones, and provided with a numeric keypad. Before the start of the experiment, subjects performed practice runs to familiarize themselves with the stimuli and task. Practice runs presented only a diotic target sequence with no masker and no spatial cues applied. Participants practiced pressing the keypad “enter” button when they heard a color word from the talker. Participants had the opportunity to discuss the task with the experimenter between practice runs. At least 3 practice runs were presented to each participant; however, each participant could perform as many practice runs as they needed to feel comfortable with the task. Participants were then instructed that there would be two streams during testing, each coming from a different location, and that the cue word ‘bag’ would be presented at the start of each block to indicate the location of the target sequence in that block.

## Behavioral Performance Estimation

In each block, we recorded the time of each button press and compared it to the onset of word tokens to determine rates of “hits” (correct responses to a target color word) and, for cases with

the speech masker, “false alarms” (incorrect responses following a color word in the distractor stream). We also defined “object responses” (response made after an object word). For each of these response types, we defined windows for scoring response rates. “Hit windows” encompassed 300-1400 ms following the onset of each color word in the target stream; “false alarm windows” spanned 300-1400 ms after each color in the masker stream; and “object response windows” covered 300-1400 ms following each object word onset in the target stream. To resolve ambiguity due to overlap of these potential windows, we removed overlapping windows with a hierarchical rule. If a potential hit window overlapped a potential false alarm window, that false alarm window was removed. Potential object response windows that overlapped with either hit or false windows were also removed. Additionally, if two potential time windows of the same type overlapped, the first of the two was kept and the subsequent window thrown out. Using this approach, all response windows were of equal length in time and did not overlap with one another. Since these windows were defined independent of when button presses occurred, they allowed us to operationally define hit rates, false alarm rates, and object response rates, in a fair and unbiased manner.

We computed hit rates (first column) and object rates for both masker types and for each color word position (lead, lag). Object rates fell near 0% for all participants and did not vary with spatialization or masker type. For the speech masker, we also calculated false alarm (FA) rates. We opted to analyze hit and FA rates because false alarms did not exist in the noise masker case, as there were no color words in the noise masker stream. In analyzing hit rates and false alarm rates separately, we also tap into two different aspects of task performance. Hit rate reflects the ability to segregate target from masker; a failure to segregate the competing streams will lead to low hit rates, because none of the words will be intelligible. In contrast, false alarm rates reflect failures to differentiate target and masker word locations. If target and masker streams are perceptually segregated and sufficiently distinct in their perceived locations, listeners should be able to focus attention on the target stream and ignore the masker stream (not responding to color words from the wrong direction).

Hit, false alarm, and object response rates were calculated as the total sum of button presses in each of the corresponding remaining time windows divided by the total number of those windows. Rates were computed individually for each subject in each condition.

## fNIRS Data Acquisition and Analysis

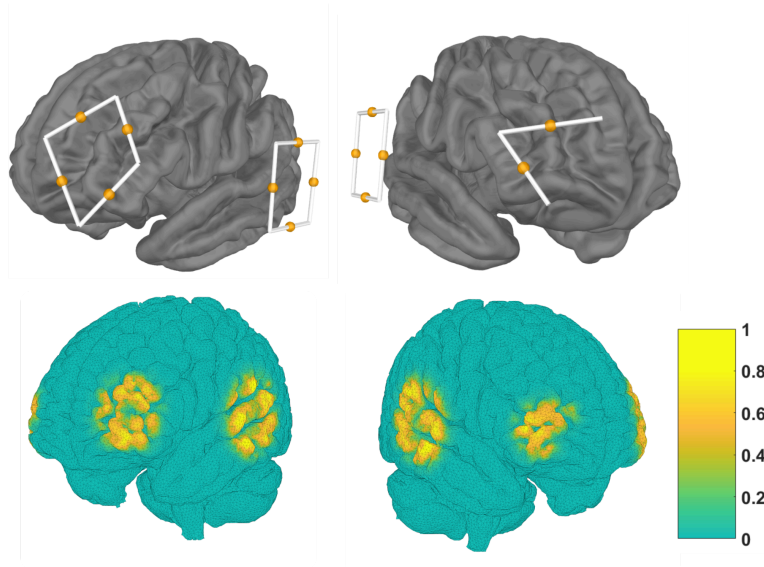
Functional near-infrared spectroscopy (fNIRS) was used to record changes in concentration of oxygenated and deoxygenated hemoglobin concentration during the task. We used a NIRx NIRSport 2 (NIRx, Berlin) system and recorded at 760 and 850 nm wavelengths with a source-detector distance of 3 cm. The source-detector montage and sensitivity map are shown in Fig. 3. Subjects wore an fNIRS head-cap fitted with light sources and detectors covering the prefrontal cortex (PFC) and superior temporal gyrus (STG) bilaterally, resulting in 14 source-detector pairs. Source and detector locations were chosen by selecting “dorsolateral prefrontal cortex” and “superior temporal gyrus” in the fNIRS Optodes’ Location Decider (fOLD)



software (Zimeo Morais et al., 2018). Short channel detectors were attached to each source in order to record systemic physiological signals from the scalp.

fNIRS data processing was performed in Python using the MNE software package (Gramfort et al., 2013; Luke et al., 2021-4). Raw intensity values were first converted to optical density. To exclude source-detector pairs (channels) with poor data quality, we calculated the scalp coupling index (SCI) for each channel, a measure of the presence of the heart rate signal in the data (Pollonini et al., 2014). Channels with an SCI below 0.8 were removed from further analysis. Out of 14 long-distance channels, across the different participants an average of  $1.42 \pm 1.47$  S.D. channels were removed due to low SCI. Next, high-frequency motion artifacts were removed using temporal derivative distribution repair (Fishburn et al., 2019). Time traces were then bandpass filtered between 0.01 and 0.3 Hz using a fourth order, zero-phase Butterworth filter to remove heartbeat and respiration signals. In order to further remove physiological artifacts, short channel information was regressed out of the optical density signal using the `short_channel_regression` function in MNE (Fabbri et al., 2004; Saager and Berger, 2005; Scholkmann et al., 2014). Finally, we applied the modified Beer-Lambert law (Kocsis et al., 2006) to the time traces of optical density to compute changes in oxygenated hemoglobin ( $\Delta\text{HbO}$ ) and deoxygenated hemoglobin ( $\Delta\text{HbR}$ ) concentrations.

The 13-18 second gap after each block ensured that hemoglobin levels returned to baseline before the next block began. The length of the silent period was chosen randomly from that range on each block in order to avoid the presence of a coherent signal in the fNIRS data at the rate of stimulus presentation. Raw fNIRS data for each block were epoched from -5 to 20 seconds from the onset of the sound cue. The first 5 seconds of each epoch (before and up to the cue) were used to set the baseline of  $\Delta\text{HbO}$  and  $\Delta\text{HbR}$  by subtracting the mean value over this period from the raw measure. Baseline was performed separately for  $\Delta\text{HbO}$  and  $\Delta\text{HbR}$  data. To quantify the strength of the hemodynamic response in each block, we calculated the mean  $\Delta\text{HbO}$  value during sound stimulation (from stimulus onset to 12.8 seconds later) in each channel for each block.



**Figure 3.** fNIRS montage and sensitivity map. Sensors covered the prefrontal cortex and superior temporal gyrus, bilaterally. The top row shows placement of near-infrared light sources and detectors (end points of white lines), and the location of data channels (between each source and detector, yellow dots). The bottom row shows the normalized sensitivity map projected onto the ICBM 2009c nonlinear asymmetric brain model (Fonov et al., 2009).

## Statistical Analysis

We used linear mixed effects regression to analyze task performance and by-channel hemodynamic response magnitudes. All models were implemented in R version 4.2.1 (R Core Team, 2024) and fit using the “lmer” function in the lme4 library (Bates et al. 2015). To estimate main effects and interactions, we used the “mixed” function in the afex library (Singmann, H., Bolker, B., Westfall, J., & Aust, F., 2018), which uses likelihood ratio tests to assess the fit of a model with the effect of interest versus that of a simplified model without that effect. Thus, main effects and interactions are reported with chi-square values that index whether inclusion of the effect leads to a significant improvement in model fit. We followed up on significant effects with Bonferroni-corrected pairwise comparisons.

To analyze task performance, we fit separate models for hit rates and false alarm rates. For both analyses, we included fixed effects of spatialization (small ITD, large ITD, natural ILD, broadband ILD) and color word position (lead, lag), as well as random intercepts for each participant. For the hit rate model only, we included a fixed effect of masker type (speech, noise). We did not include this effect for false alarms, since false alarms were not possible in the noise masker condition.

To compare hemodynamic response magnitudes, we fit separate statistical models for mean  $\Delta\text{HbO}$  in PFC, mean  $\Delta\text{HbR}$  in PFC, mean  $\Delta\text{HbO}$  in STG, and mean  $\Delta\text{HbR}$  in STG. We fit separate models for each masker type. For all fNIRS models, we included a fixed effect of

spatialization, random intercepts for each subject, and random intercepts for each channel, allowing the model to account for baseline differences in activation across channels. Finally, for the STG models only, we also included a fixed effect of cortical hemisphere (contralateral, ipsilateral), driven by our a priori hypothesis that spatial selective attention will result in an increase in neural activity in the auditory cortex contralateral to attention. We did not include a fixed effect of cortical hemisphere for the PFC data because the sensor montage was not symmetrical, and therefore did not image identical populations in left and right hemispheres.

For each analysis, the omnibus model used sum-coded contrasts for all factors (e.g., the fixed effect of cortical hemisphere might be coded with a [1,-1] contrast for [contralateral, ipsilateral]). To perform pairwise comparisons between different levels of spatialization (e.g., small ITD vs. large ITD), we re-fit the omnibus model using treatment coding. With a treatment coding scheme, one level (e.g., small ITD) was set as a reference level, and beta estimates reflect pairwise comparisons between the reference level and each other level. To obtain all pairwise estimates, we ran four models, rotating which level was the reference level. Critically, these follow-up models do not differ in their fit to the data; all that differs across the models is the comparison captured by each beta coefficient. For all pairwise comparisons, we therefore report beta coefficients that capture the difference between those conditions in the same units as the dependent measure (e.g., hit rate).

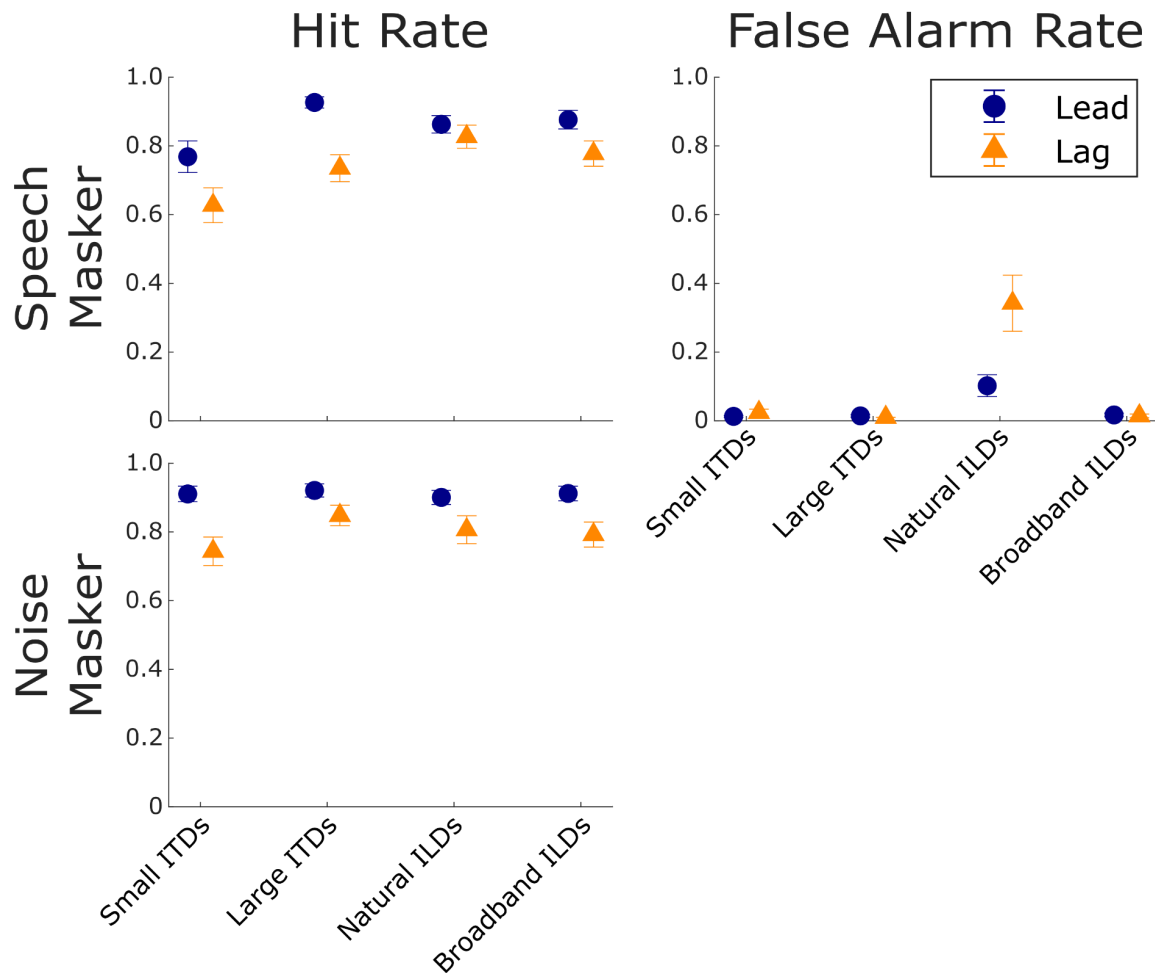
## Results

### Hit rates varied between spatialization conditions

Hit rates across spatialization, masker type, and color word position are shown in the first column of Figure 4. The hit rates for the small ITD condition (leftmost data points) tended to be lower than for any other condition. Hit rates were slightly higher for the noise masker than the speech masker (compare bottom left and top left panels). Additionally, hit rates tended to be higher for a color word that occurred in the leading position within a word token pair than for a lagging position across all spatialization conditions (compare blue circles to corresponding yellow triangles).

Statistical analyses supported these observations. The linear mixed effects model revealed a significant main effect of masker type ( $\chi^2(1) = 18.01$ ,  $p < 0.001$ ), driven by higher hit rates with the noise masker (mean: 0.854) compared to with the speech masker (mean: 0.800;  $\beta = -0.054$ ,  $p < 0.001$ ). We also observed a significant main effect of color word position ( $\chi^2(1) = 76.51$ ,  $p < 0.001$ ), driven by higher hit rates for color words occurring in the leading position (mean: 0.885) compared to in the lagging position (mean: 0.770;  $\beta = -0.115$ ,  $p < 0.001$ ). The model also revealed a significant main effect of spatialization ( $\chi^2(3) = 34.78$ ,  $p < 0.001$ ) and a significant interaction between spatialization and masker type ( $\chi^2(3) = 13.28$ ,  $p = 0.004$ ). The interaction between spatialization and position was marginally significant ( $\chi^2(3) = 6.77$ ,  $p = 0.080$ ), as was the three way interaction between spatialization, masker, and color word position ( $\chi^2(3) = 7.22$ ,  $p$

= 0.065). To interpret the significant interaction between spatialization and masker type, we conducted separate analyses of the effect of spatialization within each masker type. Within the speech masker data, the hit rates with small ITDs (mean: 0.698) were significantly lower than with large ITDs (mean: 0.830;  $\beta = 0.129$ ,  $p < 0.001$ ), natural ILDs (mean: 0.845;  $\beta = 0.147$ ,  $p < 0.001$ ), and broadband ILDs (mean: 0.827;  $\beta = 0.133$ ,  $p < 0.001$ ), but the large ITD, natural ILD, and broadband ILD hit rates did not differ significantly from one another. For the noise masker, the hit rates with small ITDs (mean: 0.864) were significantly lower than with large ITDs (mean: 0.907;  $\beta = 0.044$ ,  $p < 0.001$ ). Hit rates for all other pairs of conditions were not significantly different when the masker was noise.



**Figure 4.** Hit and False Alarm Rates. First column (Hit Rates): The percentage of responses to color words in the target stream when the masker was speech (top row) or noise (bottom row). For hit rates, significance asterisks show comparisons across spatialization, collapsed across word position (lead, lag). Second column (False Alarm Rates): The percentage of responses to color words in the masker stream only when the masker was speech. For false alarm rates, significance asterisks show comparisons across spatialization within each level of word position (lead, lag), driven by a significant interaction. For all plots, blue circles show results when a color word occurred in the leading position, and orange triangles show

results for a color word in the lagging position. Error bars indicate the mean and SEM across subjects in each spatialization condition.

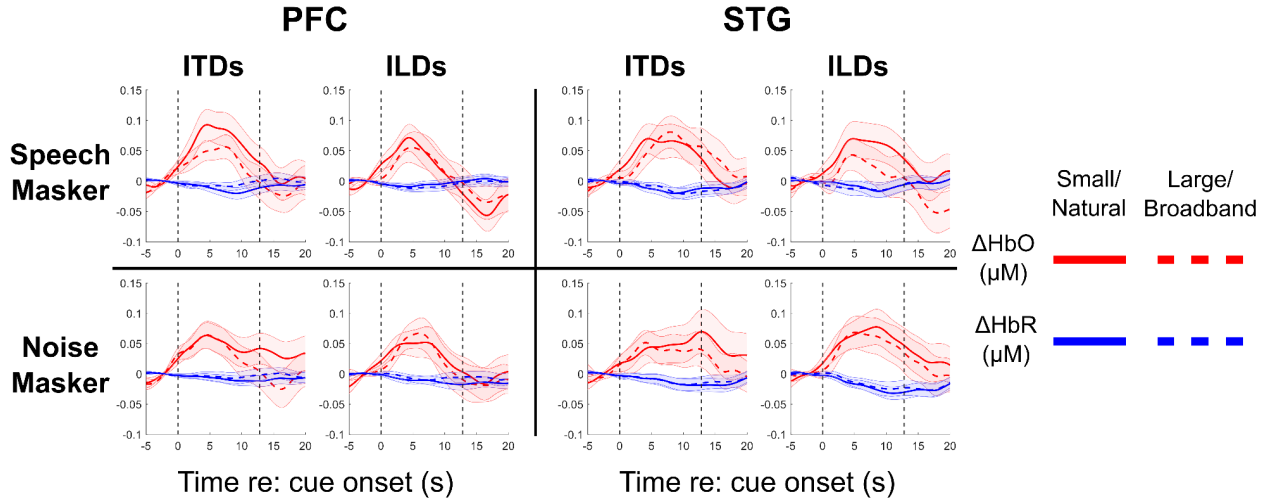
## False alarm rates were larger with natural ILDs

As noted above, false alarm rates (reflecting responses to a color word in the distracting stream) could only be computed for the speech masker conditions (top row of Figure 4). As seen in the second column of Figure 4, the false alarm rates were near zero for small ITD, large ITD, and broadband ILD conditions; however, they were higher in the natural ILD condition, especially for lagging color words.

Statistical analyses confirmed these observations: the mixed effects model showed a significant main effect of spatialization ( $\chi^2(3) = 87.42$ ,  $p < 0.001$ ), a significant main effect of color word position ( $\chi^2(1) = 12.28$ ,  $p < 0.001$ ), and a significant interaction between the two ( $\chi^2(3) = 33.47$ ,  $p < 0.001$ ). To interpret the interaction, we conducted post hoc analysis across spatialization condition within each color word position. For the leading color word, false alarm rates were significantly larger in the natural ILD condition (mean: 0.102) than in the small ITD (mean: 0.013;  $\beta = -0.089$ ,  $p < 0.001$ ), large ITD (mean: 0.014;  $\beta = -0.088$ ,  $p < 0.001$ ), and broadband ILD (mean: 0.017;  $\beta = -0.085$ ,  $p < 0.001$ ) conditions. An identical, but statistically stronger pattern arose for the lagging color word: false alarm rates were significantly larger in the natural ILD condition (mean: 0.342) than in the small ITD (mean: 0.0243;  $\beta = -0.318$ ,  $p < 0.001$ ), large ITD (mean: 0.010;  $\beta = -0.328$ ,  $p < 0.001$ ), and broadband ILD (mean: 0.014;  $\beta = -0.332$ ,  $p < 0.001$ ) conditions. None of the false alarm rates for leading and lagging color words in the small ITD, large ITD, or broadband ILD conditions differed significantly from one another.

## fNIRS Results

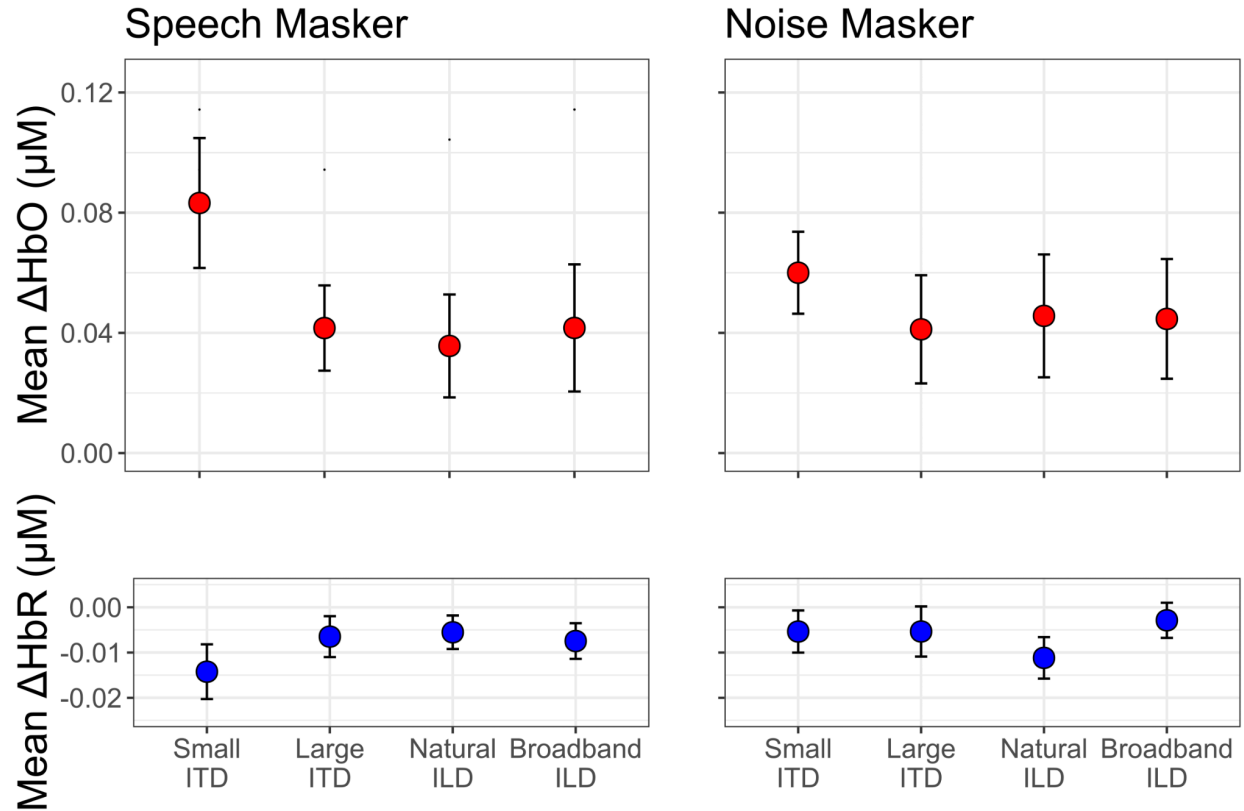
Figure 5 shows the block averaged  $\Delta\text{HbO}$  (red lines) and  $\Delta\text{HbR}$  (blue lines) mean time traces and standard errors across subjects (colored ribbons), averaged within each ROI and masker type. The left plot in each panel compares the small (solid) and large (dashed) ITD spatialization conditions, and the right plot compares natural (solid) and broadband (dashed) ILD spatializations. Hemodynamic responses in all regions and conditions show a canonical shape, with  $\Delta\text{HbO}$  increasing with respect to baseline during sound stimulation (between 0 and 12.8 seconds) and decreasing thereafter (Luke et al., 2021-4; Quaresima and Ferrari, 2019).



**Figure 5.** Block average hemodynamic responses. Each individual subplot shows the  $\Delta\text{HbO}$  (red) and  $\Delta\text{HbR}$  (blue) curves for each masker type and in each region. Top left panel: Hemodynamic responses in PFC when the masker was speech (small vs. large ITD, left subpanel; Natural vs. Broadband ILD, right subpanel). Top right panel: as before, in STG when the masker was speech. Bottom left panel: as before, in PFC when the masker was noise. Bottom right panel: as before, in STG when the masker was noise.

Figure 6 summarizes the differences in hemodynamic responses by plotting the mean  $\Delta\text{HbO}$  and  $\Delta\text{HbR}$  averaged within prefrontal cortex during the block (0-12.8 seconds from cue onset). Mean  $\Delta\text{HbO}$  for the small ITD condition was larger than for any other condition when the masker was speech, but mean  $\Delta\text{HbO}$  did not vary strongly with spatialization when the masker was noise.

Consistent with this summary, our statistical models only found a significant effect of spatialization in mean  $\Delta\text{HbO}$  in PFC ( $\chi^2(3) = 27.15$ ,  $p < 0.001$ ). Post hoc tests showed this main effect was driven by significantly larger mean  $\Delta\text{HbO}$  in the small ITD condition (mean:  $0.086 \mu\text{M}$ ) compared to the large ITD (mean:  $0.039 \mu\text{M}$ ;  $\beta = -0.047$ ,  $p < 0.001$ ), natural ILD (mean:  $0.037 \mu\text{M}$ ;  $\beta = -0.049$ ,  $p < 0.001$ ), and broadband ILD (mean:  $0.044 \mu\text{M}$ ;  $\beta = -0.042$ ,  $p < 0.001$ ) conditions, with no significant differences between the other three spatialization conditions. When the masker was noise, the effect of spatialization on mean  $\Delta\text{HbO}$  in PFC was not significant ( $\chi^2(3) = 3.82$ ,  $p = 0.282$ ). Mean  $\Delta\text{HbR}$  in PFC did not significantly vary in either the speech masker case ( $\chi^2(3) = 4.94$ ,  $p = 0.176$ ) or the noise masker case ( $\chi^2(3) = 4.48$ ,  $p = 0.214$ ).

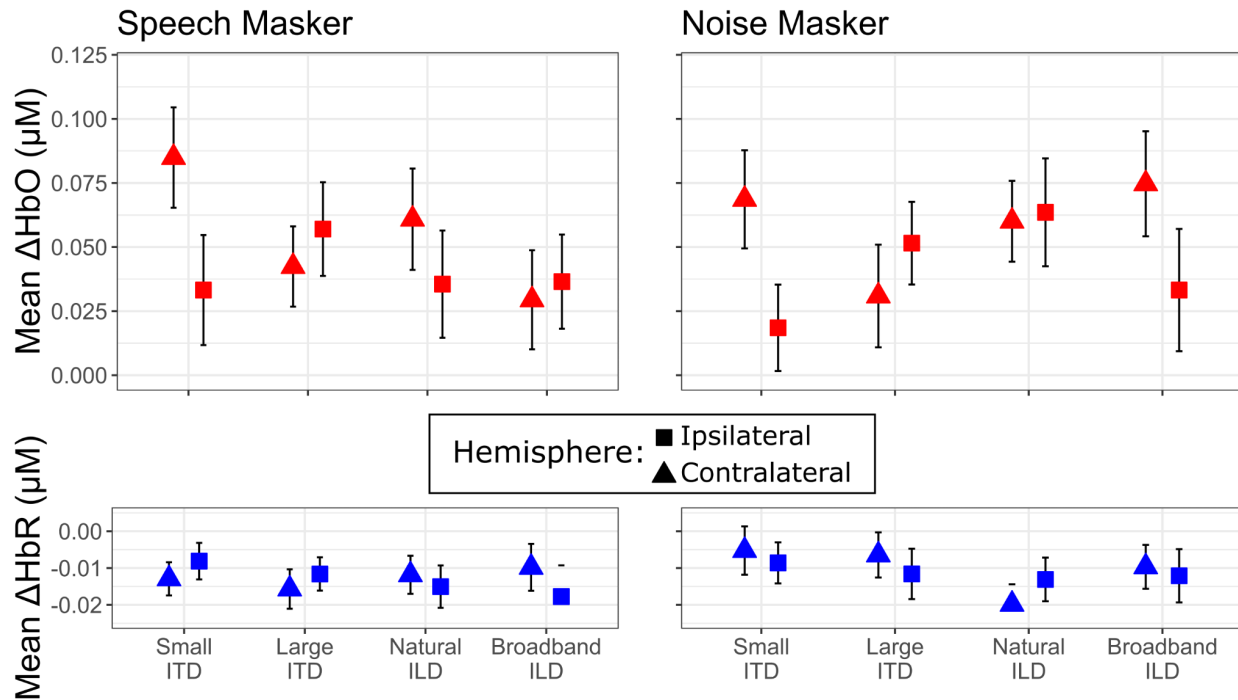


**Figure 6.** Mean hemodynamic response magnitudes in prefrontal cortex. Top row: Mean changes in oxygenated hemoglobin concentration ( $\Delta\text{HbO}$ ) when the masker was speech (left panel) and noise (right panel). Bottom row: Mean changes in deoxygenated hemoglobin concentration ( $\Delta\text{HbR}$ ) when the masker was speech (left panel) and noise (right panel).

Figure 7 shows mean  $\Delta\text{HbO}$  and mean  $\Delta\text{HbR}$  in STG averaged separately across channels in the hemisphere contralateral (triangles) and ipsilateral (squares) to the attended target direction. For the speech masker, hemodynamic responses in STG were similar across hemispheres and across spatialization cues for the large ITD, natural ILD, and broadband ILD conditions. However, in the small ITD condition, the hemodynamic response was stronger in the cortical hemisphere contralateral to the direction of attention. For the noise masker, hemodynamic responses varied more with spatial cue condition, but did not differ between hemispheres for large ITD or natural ILD conditions; however, for both small ITD and broadband ILD conditions, hemodynamic responses were larger in the hemisphere contralateral to the direction of attention.

These observations were supported by our statistical analyses. The mean  $\Delta\text{HbO}$  models showed a significant interaction between spatialization and hemisphere both when the masker was speech ( $\chi^2(3) = 7.99$ ,  $p = 0.046$ ) and when the masker was noise ( $\chi^2(3) = 11.02$ ,  $p = 0.012$ ). To follow up on this interaction, we conducted post hoc tests comparing across hemispheres within each spatialization condition. For the speech masker, the interaction effect was driven by an effect of hemisphere in the small ITD condition only, with larger mean  $\Delta\text{HbO}$  in contralateral

channels (mean: 0.078  $\mu\text{M}$ ) compared to the ipsilateral channels (mean: 0.045  $\mu\text{M}$ ;  $\beta = -0.033$ ,  $p = 0.0417$ ). A similar pattern emerged in the noise masker data, with larger responses in contralateral channels (mean: 0.0645  $\mu\text{M}$ ) than in ipsilateral channels (mean: 0.0285  $\mu\text{M}$ ;  $\beta = -0.036$ ,  $p = 0.017$ ) for the small ITD condition only; the effect of hemisphere approached significance in the broadband ILD condition with the noise masker (contralateral mean: 0.0673, ipsilateral mean: 0.0406;  $\beta = -0.027$ ,  $p = 0.077$ ).



**Figure 7.** Mean hemodynamic response magnitudes in superior temporal gyrus. Top row: Mean  $\Delta\text{HbO}$  when the masker was speech (left panel) and noise (right panel). Data within each spatialization condition are shown for sensors contralateral (triangles) and ipsilateral (squares) to the direction of the target sound stream. Bottom row: As in the top row, but for mean  $\Delta\text{HbR}$ .

## Discussion

We examined how ITDs and ILDs shaped both perceptual performance and cortical activity in a color word detection task. Performance was generally better for a noise masker than a speech masker, for leading color words than lagging color words, and for stronger spatial cues (large ITDs, broadband ILDs) versus less robust spatial differences (small ITDs and natural ILDs). When listening with a speech distractor with small ITDs, hemodynamic activity in PFC was stronger than for other conditions. For both speech and noise maskers, activity in STG was lateralized contralateral to the direction of attention when sources were lateralized with small



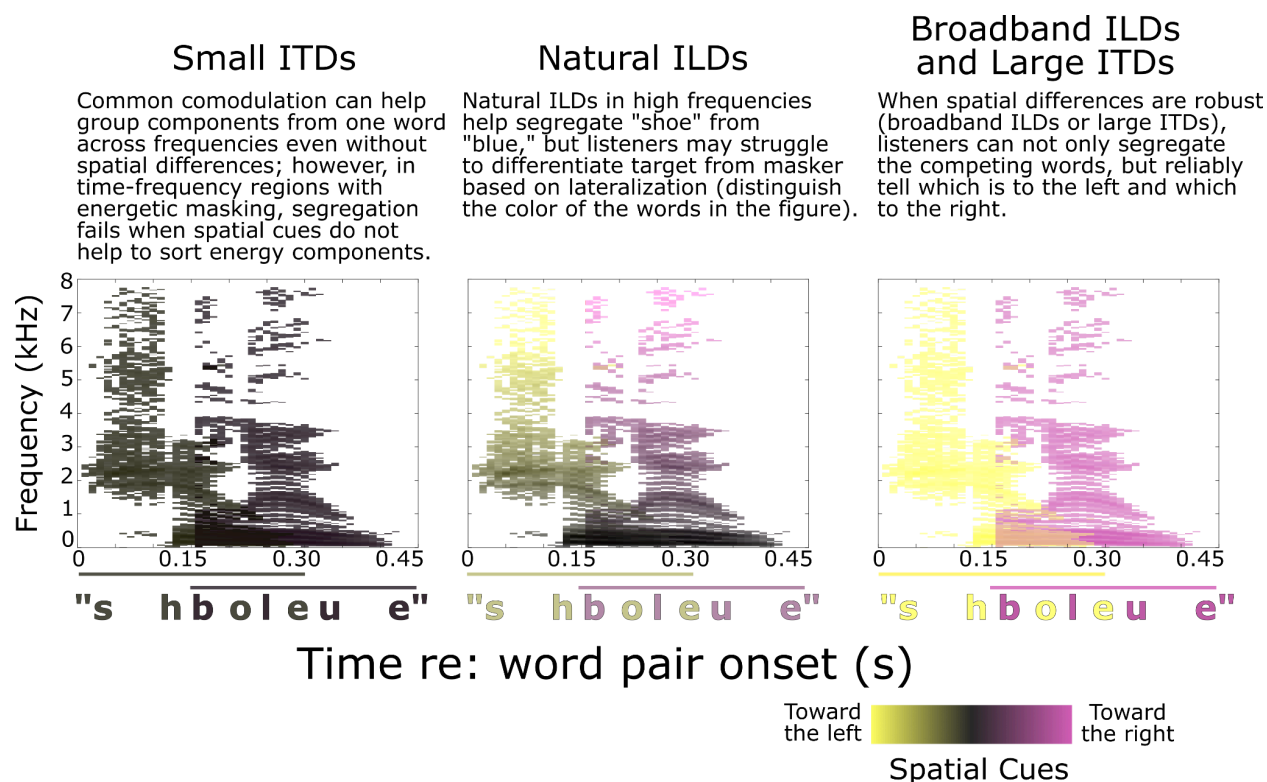
ITDs, but not for any other spatial conditions. Together, these results provide new insights into how different spatial cues contribute to SRM and speech understanding in noisy settings.

## Energetic masking accounts for noise masker results

In all spatial conditions, the hit rates were higher on average when the masker consisted of white noise bursts rather than speech. This supports the idea that participants were better able to hear out the target words when the background was spectrotemporally distinct from the target speech. This effect is well documented: Listeners perform better on SRM tasks when the masker is noise than when it is speech, even at negative signal to noise ratios (Best et al., 2012; Johnstone and Litovsky, 2006; Moore, 2008; Zhang et al., 2021a). Noise maskers do not contain any structured spectrotemporal features or semantic or linguistic content, and thus are not confused with words in the target stream. With a noise masker, both target segregation and selection are relatively trivial, making spatial information redundant.

When segregation and selection are easy (i.e., when the masker is noise), energetic masking is likely the limiting factor on performance. In these conditions, hit rates were larger for leading color words than for lagging color words. This asymmetry likely reflects the fact that the lagging token (word or noise) onset was energetically masked by the end of the leading token in each pair. Specifically, syllable onsets convey more lexical information than syllable codas (Pimentel et al., 2021; Sun and Poeppel, 2023); therefore, energetic masking of the onset of a lagging word has a greater impact on performance than masking of the coda of a leading word. For the noise masker, hit rates were lower with small ITDs than with large ITDs, which might be explained by the influence that robust ITD differences can have in combatting energetic masking (Culling and Lavandier, 2021; Durlach, 1963).

For the speech masker, hit rates were higher for leading than for lagging color words. This difference mirrors the effect for noise maskers, suggesting that energetic masking of the onsets of lagging target words degrades performance for speech maskers like it does for noise maskers. However, for speech maskers, hit rates were generally lower than for the noise masker and also varied with spatial cues, hinting that other factors affect performance, as well. These other factors are considered in the following sections.



**Figure 8.** Schematic representation of how different spatial cues impact segregation and spatial perception. Each panel shows a spectrogram representation of the word pair “shoe” (leading) and “blue” (lagging). In each time-frequency pixel, the dominant spatial cues determine the color (yellow towards the left, magenta to the right). Time-frequency pixels where there is significant overlap between words are colored by the mean of the two color schemes. The horizontal lines at the bottom of each panel denote the time extent of the lead and lag words; the color of the lines and the corresponding word labels denote the perceived laterality of the overall word, integrating spatial information across the grouped spectrum. Small ITDs may fail to support segregation, leading listeners to hear an unintelligible mixture of both words (orthographically, “shboleue,” left panel). Natural ILDs may support segregation, but the perceived words may be perceived near midline, making it hard to select the target and ignore the masker. Robust spatial cues (broadband ILDs and large ITDs) enable both segregation and spatial selection.

## ITDs and ILDs yield different modes of failure for speech maskers

In the small ITD condition, listeners had relatively low hit rates but also low false alarm rates: Listeners struggled to identify color words in either stream, suggesting that they struggled to group spectrotemporal features into intelligible words (i.e., a failure of segregation). This is schematized in the first row of Fig. 8. Without access to strong spatial cues (space is conveyed with color in Fig. 8), a listener would struggle to separate “shoe” from “blue;” instead, they may have sometimes heard a mixture of both words that was difficult to interpret (“shboleue”). This

also helps explain why the small ITD condition yielded a larger hemodynamic response in PFC than the other spatialization conditions; listeners expended greater cognitive effort in their attempt to segregate the competing words (Lawrence et al., 2018; Zhang et al., 2021b; Zhou et al., 2022a, 2022b). With small ITDs, activity in STG also lateralized, increasing in the hemisphere contralateral to the direction of attention — as is often observed when spatial auditory attention is strongly deployed (Alho et al., 1998, 1999, 2003; Ciaramitaro et al., 2007; Jäncke and Shah, 2002; Yang and Mayer, 2014).

In contrast, natural ILDs elicited both high hit rates and high false alarm rates. This pattern suggests that listeners could hear out individual word tokens, but could not differentiate whether color words belonged to the target or masker stream — a failure of selection, rather than segregation. As illustrated in the middle panel of Figure 8, ILDs in the high frequencies may help support segregation by binding with comodulated lower frequency sound components; however, the perceived lateralities of the resulting words are dominated by the lower-frequency energy, which carried zero- $\mu$ s ITDs and only modest ILDs in the natural ILD condition. Both target and masker words are thus likely to be heard near midline, leading to increased confusions about which word was from the target and which from the masker. Notably, the false alarm rate is higher for lagging than leading masker color words, suggesting that the lateralization of the lagging word is especially weak. A number of studies have shown that interaural differences at sound onsets strongly influence lateralization (Freyman and Zurek, 2017; Klein-Hennig et al., 2011). Given that the coda of the leading word energetically masks the onset of the lagging word, it is unsurprising that listeners may be more unsure of the lateralization of a lagging color masker word than a leading color masker word, leading to an even larger false alarm rate.

Although performance was relatively poor with natural ILDs, PFC activity was not stronger in this condition than with other spatializations. In the conditions in which segregation was easier (natural ILD, large ITD, and broadband ILD), PFC was less strongly engaged and STG activity less lateralized. Together, these results hint that PFC is more strongly engaged when listeners are struggling to segregate sources, but not when they are having difficulty discriminating differences in competing stream locations (for spatial selection). For small ITDs, strong engagement of spatial attention (seen in PFC) likely increases sensory responses in contralateral STG to enhance the target stream representation, albeit with limited success.

Broadband ILDs and large ITDs supported both stream segregation and spatial selection: Hit rates were high and false alarm rates were low. With large perceived spatial separation (distinct yellow and magenta in row 3 of Fig. 8), listeners both properly segregate constituent words and select the target.

These results highlight the importance of separately considering failures of segregation and spatial selection when trying to understand the processes allowing listeners to make sense of competing sounds in a mixture. Specifically, listeners must group sounds into objects based on spectrotemporal information (comodulation of different frequency components, harmonic structure, etc.). Natural ILDs alone can help support object formation of overlapping sounds by helping the brain sort out which components belong to which source. The perceived location of

each object can provide a powerful cue for guiding selective attention when sources are perceived at different locations (Allen et al., 2008; Bregman, 1990; Darwin, 2007; Kidd et al., 2010; Shinn-Cunningham, 2008b). However, even spatial cues that are sufficient to aid source segregation (like natural ILDs) may not produce robust differences in perceived location, leading to confusions between properly segregated sources.

## Summary, Limitations, and Future Directions

Here, we compared ILD and ITD spatialization conditions. We found that weaker ITDs and ILDs (small and natural, respectively) led to poorer performance than stronger ITDs and ILDs (large and broadband, respectively). Moreover, differences in the response patterns suggest that small ITDs led to failures of segregation, while natural ILDs, which are modest in the frequency range where most speech energy resides, produced more selection errors. It is tempting to conclude that ITDs and ILDs thus contribute differently to SRM performance, in general. However, we did not match these “weak” ITD and ILD conditions to produce perceptually similar lateralization of our stimuli, limiting our ability to directly compare across conditions.

In the small ITD condition, we observed greater hemodynamic responses in PFC and a lateralization of hemodynamic responses in STG, with greater activation in the hemisphere contralateral to the direction of the target. Considered alongside our behavioral results, this pattern suggests that PFC was more strongly engaged in the one condition where segregation failures limited performance. Although these findings are intriguing, they are not definitive. The measures of hemodynamic activity obtained with this fNIRS montage were relatively coarse with four channels of data per hemisphere of STG and seven channels of data in PFC. We also gathered only eight blocks of data per condition. Further investigations employing denser fNIRS montages, presenting more blocks in each spatialization condition, and co-registering results with MRI structural scans would provide more definitive results (Yücel et al., 2024). Such studies would provide better ground-truth estimates of the neural locus of activity with better signal-to-noise.

Future experiments should investigate psychophysically matched ITDs and ILDs or parametrically vary ITD and ILD magnitudes to support direct comparisons between ITDs and ILDs that produce similar spatial percepts. Clinically focused experiments could also explore whether improved spatial hearing outcomes in bilateral cochlear implant users with magnified ILDs (Brown, 2014, 2018; Richardson et al., 2025) are driven by the same mechanisms as in normal hearing listeners.

## Author Contributions

Conceptualization: Benjamin N. Richardson, Barbara G. Shinn-Cunningham, Christopher A. Brown

Formal analysis: Benjamin Richardson

Funding acquisition: Barbara G. Shinn-Cunningham, Christopher A. Brown

Methodology: Benjamin N. Richardson, Sahil Luthra, Jana M. Kainerstorfer, Barbara G. Shinn-Cunningham, Christopher A. Brown

Data collection: Benjamin N. Richardson

Resources: Jana M. Kainerstorfer, Barbara G. Shinn-Cunningham, Christopher A. Brown

Supervision: Jana M. Kainerstorfer, Barbara G. Shinn-Cunningham, Christopher A. Brown

Visualization: Benjamin N. Richardson

Writing – original draft: Benjamin N. Richardson

Writing – review & editing: Benjamin N. Richardson, Sahil Luthra, Jana M. Kainerstorfer, Barbara G. Shinn-Cunningham, Christopher A. Brown

## Data Availability Statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

## Ethics Statement

This study was reviewed and approved by The University of Pittsburgh Institutional Review Board. The patients/participants provided their written informed consent to participate prior to their participation in this study.

## Acknowledgements

The authors would like to thank Sydney Sepkovic, Cece Restauri, Kelsey Overbay, Ryley Watt, Katherine Bart, and Maanasa Guru Adimurthy for their assistance in data collection. The authors would also like to thank Eli Bulger and Victoria Figarola for their insights on fNIRS data analysis and writing preparation.

# References

- Alho, K., Connolly, J. F., Cheour, M., Lehtokoski, A., Huottilainen, M., Virtanen, J., Aulanko, R., et al. (1998). "Hemispheric lateralization in preattentive processing of speech sounds," *Neurosci. Lett.*, **258**, 9–12. doi:10.1016/s0304-3940(98)00836-2
- Alho, K., Medvedev, S. V., Pakhomov, S. V., Roudas, M. S., Tervaniemi, M., Reinikainen, K., Zeffiro, T., et al. (1999). "Selective tuning of the left and right auditory cortices during spatially directed attention," *Brain Res. Cogn. Brain Res.*, **7**, 335–341. doi:10.1016/s0926-6410(98)00036-6
- Alho, K., Vorobyev, V. A., Medvedev, S. V., Pakhomov, S. V., Roudas, M. S., Tervaniemi, M., van Zuijen, T., et al. (2003). "Hemispheric lateralization of cerebral blood-flow changes during selective listening to dichotically presented continuous speech," *Brain Res. Cogn. Brain Res.*, **17**, 201–211. doi:10.1016/s0926-6410(03)00091-0
- Allen, K., Carlile, S., and Alais, D. (2008). "Contributions of talker characteristics and spatial location to auditory streaming," *J. Acoust. Soc. Am.*, **123**, 1562–1570. doi:10.1121/1.2831774
- Arbogast, T. L., Mason, C. R., and Kidd, G., Jr (2002). "The effect of spatial separation on informational and energetic masking of speech," *J. Acoust. Soc. Am.*, **112**, 2086–2098. doi:10.1121/1.1510141
- Best, V., Gallun, F. J., Carlile, S., and Shinn-Cunningham, B. G. (2007). "Binaural interference and auditory grouping," *J. Acoust. Soc. Am.*, **121**, 1070–1076. doi:10.1121/1.2407738
- Best, V., Marrone, N., Mason, C. R., and Kidd, G., Jr (2012). "The influence of non-spatial factors on measures of spatial release from masking," *J. Acoust. Soc. Am.*, **131**, 3103–3110. doi:10.1121/1.3693656
- Bregman, A. S. (1990). *Auditory Scene Analysis: The perceptual organization of sound*, The

- MIT Press. doi:10.7551/mitpress/1486.001.0001
- Brown, C. A. (2014). "Binaural enhancement for bilateral cochlear implant users," *Ear Hear.*, **35**, 580–584. doi:10.1097/AUD.0000000000000044
- Brown, C. A. (2018). "Corrective binaural processing for bilateral cochlear implant patients," *PLoS One*, **13**, e0187965. doi:10.1371/journal.pone.0187965
- Cherry, E. C. (1953). "Some experiments on the recognition of speech, with one and with two ears," *J. Acoust. Soc. Am.*, **25**, 975–979. doi:10.1121/1.1907229
- Ciaramitaro, V. M., Buracas, G. T., and Boynton, G. M. (2007). "Spatial and cross-modal attention alter responses to unattended sensory information in early visual and auditory human cortex," *J. Neurophysiol.*, **98**, 2399–2413. doi:10.1152/jn.00580.2007
- Culling, J. F., and Lavandier, M. (2021). "Binaural unmasking and spatial release from masking," *Springer Handbook of Auditory Research*, Springer Handbook of Auditory Research, Springer International Publishing, Cham, pp. 209–241. doi:10.1007/978-3-030-57100-9\_8
- Culling, J. F., and Summerfield, Q. (1995). "Perceptual separation of concurrent speech sounds: absence of across-frequency grouping by common interaural delay," *J. Acoust. Soc. Am.*, **98**, 785–797. doi:10.1121/1.413571
- Darwin, C. J. (1997). "Auditory grouping," *Trends Cogn. Sci.*, **1**, 327–333. doi:10.1016/S1364-6613(97)01097-8
- Darwin, C. J. (2007). "Spatial Hearing and Perceiving Sources," *Auditory Perception of Sound Sources*, Springer US, Boston, MA, pp. 215–232. doi:10.1007/978-0-387-71305-2\_8
- Dorman, M. F., Loiselle, L., Stohl, J., Yost, W. A., Spahr, A., Brown, C., and Cook, S. (2014). "Interaural level differences and sound source localization for bilateral cochlear implant patients," *Ear Hear.*, **35**, 633–640. doi:10.1097/AUD.0000000000000057
- Durlach, N. I. (1963). "Equalization and cancellation theory of binaural masking-level differences," *J. Acoust. Soc. Am.*, **35**, 1206–1218. doi:10.1121/1.1918675
- Ellinger, R. L., Jakien, K. M., and Gallun, F. J. (2017). "The role of interaural differences on

- speech intelligibility in complex multi-talker environments,” *J. Acoust. Soc. Am.*, **141**, EL170. doi:10.1121/1.4976113
- Fabbri, F., Sassaroli, A., Henry, M. E., and Fantini, S. (2004). “Optical measurements of absorption changes in two-layered diffusive media,” *Phys. Med. Biol.*, **49**, 1183–1201. doi:10.1088/0031-9155/49/7/007
- Fishburn, F. A., Ludlum, R. S., Vaidya, C. J., and Medvedev, A. V. (2019). “Temporal Derivative Distribution Repair (TDDR): A motion correction method for fNIRS,” *Neuroimage*, **184**, 171–179. doi:10.1016/j.neuroimage.2018.09.025
- Fonov, V. S., Evans, A. C., McKinstry, R. C., Almlí, C. R., and Collins, D. L. (2009). “Unbiased nonlinear average age-appropriate brain templates from birth to adulthood,” *Neuroimage*, **47**, S102. doi:10.1016/s1053-8119(09)70884-5
- Freyman, R. L., Helfer, K. S., McCall, D. D., and Clifton, R. K. (1999). “The role of perceived spatial separation in the unmasking of speech,” *J. Acoust. Soc. Am.*, **106**, 3578–3588. doi:10.1121/1.428211
- Freyman, R. L., and Zurek, P. M. (2017). “Strength of onset and ongoing cues in judgments of lateral position,” *J. Acoust. Soc. Am.*, **142**, 206. doi:10.1121/1.4990020
- Gardner, W. G., and Martin, K. D. (1995). “HRTF measurements of a KEMAR,” *J. Acoust. Soc. Am.*, **97**, 3907–3908. doi:10.1121/1.412407
- Glyde, H., Buchholz, J. M., Dillon, H., Cameron, S., and Hickson, L. (2013). “The importance of interaural time differences and level differences in spatial release from masking,” *J. Acoust. Soc. Am.*, **134**, EL147–EL152. doi:10.1121/1.4812441
- Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Goj, R., et al. (2013). “MEG and EEG data analysis with MNE-Python,” *Front. Neurosci.*, Retrieved from <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2013.00267>. Retrieved from <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2013.00267>



- Gray, W. O., Mayo, P. G., Goupell, M. J., and Brown, A. D. (2021). "Transmission of Binaural Cues by Bilateral Cochlear Implants: Examining the Impacts of Bilaterally Independent Spectral Peak-Picking, Pulse Timing, and Compression," *Trends in Hearing*, **25**, 23312165211030411. doi:10.1177/23312165211030411
- Higgins, N. C., McLaughlin, S. A., Rinne, T., and Stecker, G. C. (2017). "Evidence for cue-independent spatial representation in the human auditory cortex during active listening," *Proceedings of the National Academy of Sciences*, **114**, E7602–E7611. doi:10.1073/pnas.1707522114
- Ihlefeld, A., and Shinn-Cunningham, B. (2008). "Disentangling the effects of spatial cues on selection and formation of auditory objects," *J. Acoust. Soc. Am.*, **124**, 2224–2235. doi:10.1121/1.2973185
- Jäncke, L., and Shah, N. J. (2002). "Does dichotic listening probe temporal lobe functions?," *Neurology*, **58**, 736–743. doi:10.1212/wnl.58.5.736
- Johnstone, P. M., and Litovsky, R. Y. (2006). "Effect of masker type and age on speech intelligibility and spatial release from masking in children and adults," *J. Acoust. Soc. Am.*, **120**, 2177–2189. doi:10.1121/1.2225416
- Kidd, G., Jr, Best, V., and Mason, C. R. (2008). "Listening to every other word: examining the strength of linkage variables in forming streams of speech," *J. Acoust. Soc. Am.*, **124**, 3793–3802. doi:10.1121/1.2998980
- Kidd, G., Jr, Mason, C. R., Best, V., and Marrone, N. (2010). "Stimulus factors influencing spatial release from speech-on-speech masking," *J. Acoust. Soc. Am.*, **128**, 1965–1978. doi:10.1121/1.3478781
- Klein-Hennig, M., Dietz, M., Hohmann, V., and Ewert, S. D. (2011). "The influence of different segments of the ongoing envelope on sensitivity to interaural time delays," *J. Acoust. Soc. Am.*, **129**, 3856–3872. doi:10.1121/1.3585847
- Kocsis, L., Herman, P., and Eke, A. (2006). "The modified Beer-Lambert law revisited," *Phys.*

- Med. Biol., **51**, N91–8. doi:10.1088/0031-9155/51/5/N02
- Kong, L., Michalka, S. W., Rosen, M. L., Sheremata, S. L., Swisher, J. D., Shinn-Cunningham, B. G., and Somers, D. C. (2014). "Auditory spatial attention representations in the human cerebral cortex," *Cereb. Cortex*, **24**, 773–784. doi:10.1093/cercor/bhs359
- Kubovy, M. (1988). "Should we resist the seductiveness of the space:time::vision:audition analogy?," *J. Exp. Psychol. Hum. Percept. Perform.*, **14**, 318–320.  
doi:10.1037/0096-1523.14.2.318
- Kubovy, M., and Van Valkenburg, D. (2001). "Auditory and visual objects," *Cognition*, **80**, 97–126. doi:10.1016/s0010-0277(00)00155-4
- Lawrence, R. J., Wiggins, I. M., Anderson, C. A., Davies-Thompson, J., and Hartley, D. E. H. (2018). "Cortical correlates of speech intelligibility measured using functional near-infrared spectroscopy (fNIRS)," *Hear. Res.*, **370**, 53–64. doi:10.1016/j.heares.2018.09.005
- Luke, R., Larson, E., Shader, M. J., Innes-Brown, H., Van Yper, L., Lee, A. K. C., Sowman, P. F., et al. (2021-4). "Analysis methods for measuring passive auditory fNIRS responses generated by a block-design paradigm," *Neurophotronics*, **8**, 025008.  
doi:10.1117/1.NPh.8.2.025008
- Michalka, S. W., Kong, L., Rosen, M. L., Shinn-Cunningham, B. G., and Somers, D. C. (2015). "Short-Term Memory for Space and Time Flexibly Recruit Complementary Sensory-Biased Frontal Lobe Attention Networks," *Neuron*, **87**, 882–892. doi:10.1016/j.neuron.2015.07.028
- Middlebrooks, J. C. (2017). "Spatial Stream Segregation," In J. C. Middlebrooks, J. Z. Simon, A. N. Popper, and R. R. Fay (Eds.), *The Auditory System at the Cocktail Party*, Springer International Publishing, Cham, pp. 137–168. doi:10.1007/978-3-319-51662-2\_6
- Middlebrooks, J. C., and Waters, M. F. (2020). "Spatial Mechanisms for Segregation of Competing Sounds, and a Breakdown in Spatial Hearing," *Front. Neurosci.*, **14**, 571095.  
doi:10.3389/fnins.2020.571095
- Moore, B. C. J. (2008). "Basic auditory processes involved in the analysis of speech sounds,"

- Philos. Trans. R. Soc. Lond. B Biol. Sci., **363**, 947–963. doi:10.1098/rstb.2007.2152
- Moore, B. C. J., and Gockel, H. (2002). “Factors Influencing Sequential Stream Segregation,” *Acta Acustica united with Acustica*, **88**, 320–333. Retrieved from <https://www.ingentaconnect.com/content/dav/aaua/2002/00000088/00000003/art00004>
- Noyce, A. L., Kwasa, J. A. C., and Shinn-Cunningham, B. G. (2023). “Defining attention from an auditory perspective,” *Wiley Interdiscip. Rev. Cogn. Sci.*, **14**, e1610. doi:10.1002/wcs.1610
- Noyce, A. L., Lefco, R. W., Brissenden, J. A., Tobyne, S. M., Shinn-Cunningham, B. G., and Somers, D. C. (2022). “Extended Frontal Networks for Visual and Auditory Working Memory,” *Cereb. Cortex*, **32**, 855–869. doi:10.1093/cercor/bhab249
- Pimentel, T., Cotterell, R., and Roark, B. (2021). *Disambiguatory Signals are Stronger in Word-initial Positions* arXiv [cs.CL],. doi:10.48550/arXiv.2102.02183
- Pollonini, L., Olds, C., Abaya, H., Bortfeld, H., Beauchamp, M. S., and Oghalai, J. S. (2014). “Auditory cortex activation to natural speech and simulated cochlear implant speech measured with functional near-infrared spectroscopy,” *Hear. Res.*, **309**, 84–93. doi:10.1016/j.heares.2013.11.007
- Quaresima, V., and Ferrari, M. (2019). “Functional Near-Infrared Spectroscopy (fNIRS) for Assessing Cerebral Cortex Function During Human Behavior in Natural/Social Situations: A Concise Review,” *Organizational Research Methods*, **22**, 46–68. doi:10.1177/1094428116658959
- R Core Team (2024). “R: A Language and Environment for Statistical Computing,” *MSOR connections*, **1**, 23–25. Retrieved from <https://www.R-project.org/>.
- Richardson, B. N., Kainerstorfer, J. M., Shinn-Cunningham, B. G., and Brown, C. A. (2025). “Magnified interaural level differences enhance binaural unmasking in bilateral cochlear implant users,” *J. Acoust. Soc. Am.*, **157**, 1045–1056. doi:10.1121/10.0034869
- Saager, R. B., and Berger, A. J. (2005). “Direct characterization and removal of interfering absorption trends in two-layer turbid media,” *J. Opt. Soc. Am. A Opt. Image Sci. Vis.*, **22**,

1874–1882. doi:10.1364/josaa.22.001874

Scholkmann, F., Metz, A. J., and Wolf, M. (2014). “Measuring tissue hemodynamics and oxygenation by continuous-wave functional near-infrared spectroscopy--how robust are the different calculation methods against movement artifacts?,” *Physiol. Meas.*, **35**, 717–734. doi:10.1088/0967-3334/35/4/717

Shinn-Cunningham, B. G. (2008a). “Object-based auditory and visual attention,” *Trends Cogn. Sci.*, **12**, 182–186. doi:10.1016/j.tics.2008.02.003

Shinn-Cunningham, B. G. (2008b). “Object-based auditory and visual attention,” *Trends Cogn. Sci.*, **12**, 182–186. doi:10.1016/j.tics.2008.02.003

Singmann, H., Bolker, B., Westfall, J., & Aust, F. (2018). “afex: Analysis of Factorial Experiments. R package version 0.21-2,” Retrieved from <https://CRAN.R-project.org/package=afex>. Retrieved from <https://CRAN.R-project.org/package=afex>.

*SoundStretch Audio Processing Utility*, (n.d.). Available:

<https://www.surina.net/soundtouch/soundstretch.html>, (date last viewed: 21-Mar-25).

Retrieved March 21, 2025, from <https://www.surina.net/soundtouch/soundstretch.html>

Sun, Y., and Poeppel, D. (2023). “Syllables and their beginnings have a special role in the mental lexicon,” *Proc. Natl. Acad. Sci. U. S. A.*, **120**, e2215710120. doi:10.1073/pnas.2215710120

Van Rossum, G., and Drake, F. L., Jr (2009). *Python 3 Reference Manual: (Python Documentation Manual Part 2)*, Createspace, 244 pages. Retrieved from [https://books.google.com/books/about/Python\\_3\\_Reference\\_Manual.html?id=KlybQQAACAAJ](https://books.google.com/books/about/Python_3_Reference_Manual.html?id=KlybQQAACAAJ)

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., et al. (2020). “SciPy 1.0: fundamental algorithms for scientific computing in Python,” *Nat. Methods*, **17**, 261–272. doi:10.1038/s41592-019-0686-2

- Wightman, F. L., and Kistler, D. J. (1992). "The dominant role of low-frequency interaural time differences in sound localization," J. Acoust. Soc. Am., **91**, 1648–1661.  
doi:10.1121/1.402445
- Yang, Z., and Mayer, A. R. (2014). "An event-related fMRI study of exogenous orienting across vision and audition: Cross-Modal Orienting," Hum. Brain Mapp., **35**, 964–974.  
doi:10.1002/hbm.22227
- Yost, W. A., and Dye, R. H., Jr (1988). "Discrimination of interaural differences of level as a function of frequency," J. Acoust. Soc. Am., **83**, 1846–1851. doi:10.1121/1.396520
- Yücel, M. A., Luke, R., Mesquita, R. C., von Lühmann, A., Mehler, D. M. A., Lühns, M., Gemignani, J., et al. (2024). *The fNIRS reproducibility study hub (FRESH): Exploring variability and enhancing transparency in fNIRS neuroimaging research* BITSS,.  
doi:10.31222/osf.io/pc6x8
- Zhang, M., Alamatsaz, N., and Ihlefeld, A. (2021a). "Hemodynamic Responses Link Individual Differences in Informational Masking to the Vicinity of Superior Temporal Gyrus," Front. Neurosci., **15**, 675326. doi:10.3389/fnins.2021.675326
- Zhang, M., Alamatsaz, N., and Ihlefeld, A. (2021b). *Hemodynamic responses link individual differences in informational masking to the vicinity of superior temporal gyrus p.* 2020.08.21.261222. Retrieved from  
<https://www.biorxiv.org/content/10.1101/2020.08.21.261222v2>
- Zhou, X., Burg, E., Kan, A., and Litovsky, R. Y. (2022a). "Investigating effortful speech perception using fNIRS and pupillometry measures," Current Research in Neurobiology, **3**, 100052. doi:10.1016/j.crneur.2022.100052
- Zhou, X., Sobczak, G. S., McKay, C. M., and Litovsky, R. Y. (2022b). "Effects of degraded speech processing and binaural unmasking investigated using functional near-infrared spectroscopy (fNIRS)," PLoS One, **17**, e0267588. doi:10.1371/journal.pone.0267588
- Zimeo Morais, G. A., Balardin, J. B., and Sato, J. R. (2018). "fNIRS Optodes' Location Decider

(fOLD): a toolbox for probe arrangement guided by brain regions-of-interest,” *Sci. Rep.*, **8**, 3341. doi:10.1038/s41598-018-21716-z