

Cocktail Party Problem & Auditory Objects

How is a complex auditory scene consisting of multiple auditory objects/ streams represented in human auditory cortex? What is the neural correlate of the segregation of the auditory scene into auditory objects or streams?

Is each auditory object represented by a distinct neural code? If so, when and where does this auditory scene analysis process occur? Is it robust against bottom-up acoustic saliency (e.g. the loudness of each object) and acoustic/perceptual similarity between the auditory objects?

We address these questions by recording from human subjects who selectively listen to simultaneous natural auditory narrations, using the noninvasive physiological method of magnetoencephalography (MEG).

Methods

Stimulus & Procedures: Two speakers were mixed into a single acoustic channel presented diotically. Listeners were instructed to attend to one speaker and answer comprehension questions after every 1 minute section (2 sections per condition). The listeners switched attention to the other speaker when the same stimulus repeated. All was repeated 3 times, resulting in 3 trials in each attentional condition.

In the *main experiment* (N=11), the two speakers were of different gender and were mixed with equal intensity. In the varying target-to-masker ratio (TMR) experiment (N=6), the same two speakers were used, but the intensity of one speaker was varied. In the *same gender experiment (N=3)*, both speakers were female, and the listeners underwent a training session. **MEG:** 157-channel, whole-head MEG. 1 kHz sampling rate, downsampled to 40 Hz. The neural source of MEG activity is localized using an equivalent current dipole model, one per hemisphere.

Neural Reconstruction: We reconstructed the envelope of speech using a linear decoder that integrates MEG activity over time and sensors. MEG Responses Decoder Stimulus Envelope



STRF: The spectro-temporal response function (STRF) models the neural response evoked by a unit power increase in the stimulus (by frequency). MEG Response STRF Stimulus

STRFs were estimated by boosting. Since two speakers were presented, the full model is: Response = Speaker₁ * STRF₁ + Speaker₂ * STRF₂.

Conclusions:

A complex auditory scene containing two simultaneous speakers mixed into a single acoustic channel is behaviorally, and neurally, *parsed into* auditory objects, in auditory cortex.

Attention routes the neural representation of auditory object, i.e. speaker, into distinct spatialtemporal neural networks. The attended object is more strongly represented in posterior association auditory cortex at a latency of ~100 ms. Supported by NIH R01 DC-005660

The Neural Encoding of Auditory Objects while Listening to Competing Speakers

Nai Ding¹, Jonathan Z. Simon^{1,2}

¹Department of Electrical & Computer Engineering, ²Department of Biology University of Maryland College Park







Ding and Simon (2012) Emergence of neural encoding of auditory objects while listening to competing speakers, PNAS 109, 11854-11859.

response to the speech mixture.

Different envelopes (in the upper and lower panels) are decoded from stimulus, depending on whether the listener attends to one or the other

The grand averaged correlation envelope and actual envelope of the stimulus is shown in the left.

Two decoders (spatial-temporal weighting matrices) are designed to reconstruct the attended and unattended speaker respectively.

These two decoders integrate spatial-temporal neural activity



The M100_{STRF} is significantly modulated by attention while the shorter latency response (component M50_{STRF} is not. Neither response peak is affected by the intensity change of the two speakers.

The neural source of the M100_{STRF} is roughly consistent with that of the M100 evoked by a tone pip, which is commonly localized to planum temporale (PT). The neural source of the M50_{STRF} is more anterior than the neural sources of the M100_{STRF} and M100, and therefore is more consistent with core auditory cortex.

Cortical representations are transformed from feature-based to object-based up the cortical hierarchy: from shorter latency activity in core auditory cortex to longer latency activity in posterior association auditory cortex.

Computational Sensorimotor Systems Lab



Speaker Intensity Invariant Representation

To test how robust the speakerspecific representation is, the intensity ratio between the two speakers is varied between -8 dB

The envelope of the attended speaker can be reliably decoded at all test TMRs, and the decoding performance is TMR independent.

Subjectively rated intelligibility decreases with TMR, though not the percent of questions correctly answered ($\sim 70\%$).

The TMR-independent responses suggest that a speaker-specific neural adaptation to sound intensity compensates variations in speaker intensity.

Segregation of Speakers of the Same Gender

The attended and unattended speakers are differentially represented even when the two speakers are of the same gender.

(74% questions correctly answered)

The speaker-specific cortical representations are robustly formed

Encoding Model Demonstrates Emergence

Temporal Profile of the STRF in the Varying-Loudness Experiment Influence of Relative Intensity **Attentional Modulation**





