1	
2	
3	
4	The elicitation of audiovisual steady-state responses: multi-
5	sensory signal congruity and phase effects
6	
7	Julian Jenkins III ^{1,3} , Ariane E. Rhone ^{2,3} , William J. Idsardi ^{2,3} ,
8	and David Poeppel ^⁴
9	
10	 Department of Biology, University of Maryland, College Park
11	2. Department of Linguistics, University of Maryland, College Park
12	3. Cognitive Neuroscience of Language Laboratory, University of Maryland,
13	College Park
14	4. Department of Psychology, New York University
15	
10	
18	
19	
20	
21 22	Address for correspondence:
22 23	Address for correspondence.
24	Julian Jenkins III,M.S.
25	Department of Biology
26	University of Maryland, College Park
27	College Park, MD, USA 20742
20 29	Julianj@ullu.euu
30	
31	
32	Running head: Bimodal SSR in MEG
33	

1 Abstract

2

3 Most ecologically natural sensory experiences are not limited to a single modality. 4 While it is possible to use real ecological materials as experimental stimuli, parametric 5 control of such tokens is limited. By using artificial bimodal stimuli composed of 6 approximations to ecological signals, it can be possible to observe the interactions 7 between putatively relevant stimulus attributes. Here we use MEG as an 8 electrophysiological tool and employ as a measure the steady-state response (SSR), an 9 experimental paradigm typically applied to unimodal signals. We quantify the responses 10 to a bimodal audio-visual signal with different degrees of temporal (phase) congruity, 11 focusing on properties critical to audiovisual speech. An amplitude modulated auditory 12 signal ('pseudo-envelope') is paired with a radius-modulated disc ('pseudo-mouth'), with 13 the low-frequency modulations occurring in phase or at offset phase values. We 14 observe (i) that it is possible to elicit an SSR to bimodal signals; (ii) that bimodal signals 15 exhibit greater response power than unimodal signals; and iii) that the SSR power 16 differentially reflects the congruity between signal components. The experimental 17 paradigm facilitates a quantitative characterization of properties of multi-sensory speech 18 and other bimodal computations. 19 20 Keywords: audio-visual, cross-modal, magnetoencephalography, speech, multi-sensory

21

2 Introduction

3

4 The majority of sensory experiences are not limited to a single modality and thus require 5 the observer to not only segregate information into separate objects or streams but also 6 to integrate related information into a coherent percept across sensory modalities as 7 well as across space and time (Amedi A et al., 2005; Kelly SP et al., 2008; Lalor EC et 8 al., 2007; Macaluso E and J Driver, 2005; Miller BT and M D'Esposito, 2005; Molholm S et al., 2007; Molholm S et al., 2004; Molholm S et al., 2002; Murray MM et al., 2005; 9 10 Senkowski D et al., 2006). The ability to integrate information not only unifies the 11 perception of events, but the presence of redundant information also facilitates 12 recognition, increases signal-to-noise ratio and decreases reaction times to cross-modal 13 events (Driver J and C Spence, 1998; Hershenson M, 1962; Senkowski D et al., 2006; 14 Stein BE et al., 1989). Studies examining the simultaneous serial and parallel 15 computations and physiological responses underlying the integration of information and 16 the cognition of a unified percept have important implications for advancing the 17 understanding of the binding of cross-modal information for ecologically valid behaviors 18 such as motion perception and speech recognition and comprehension (Baumann O 19 and MW Greenlee, 2007; Lakatos P et al., 2008; Miller BT and M D'Esposito, 2005; 20 Schroeder CE and P Lakatos, 2009; Schroeder CE et al., 2008) 21 22 While it has traditionally been thought that processing of cross-modal events occurs

23 primarily in association cortices (Jones EG and TP Powell, 1970; Mesulam MM, 1998),

24 recent evidence indicates that information from other sensory modalities can influence

25 cortical areas conventionally assumed to be unimodal. Electroencephalographic (EEG), 26 functional magnetic resonance (fMRI) and magnetoencephalographic (MEG) studies in 27 humans have provided evidence that visual and somatosensory signals can influence 28 neuronal activity in the auditory cortex (e.g., see Schroeder & Foxe (Schroeder CE and 29 J Foxe, 2005) for a review). Intracranial recordings and anatomical tracings in 30 macagues have affirmed the existence of multisensory inputs to unimodal cortical areas 31 (Kayser C et al., 2008). In humans, several functional imaging and intracranial studies have identified cortical networks involved in object recognition, auditory-somatosensory 32 33 and visual-somatosensory processing and integration of audio-visual speech (Calvert 34 GA et al., 1999; Calvert GA et al., 2000; Calvert GA et al., 2001; Molholm S et al., 2004; 35 Molholm S et al., 2006; Senkowski D et al., 2008). Human imaging studies have 36 identified the superior colliculus, superior temporal sulcus, intraparietal sulcus, insula 37 and several frontal cortical areas as being involved in crossmodal computation (Calvert 38 GA et al., 2001). With regard to speech, the traditional speech areas (perisylvian) have 39 been implicated as well as superior parietal, inferior parietal, inferior frontal, superior 40 temporal sulcus and left claustrum (Calvert GA et al., 2000; Campbell R, 2008; Fort A et 41 al., 2002; Olson IR et al., 2002). These findings also emphasize the importance of rapid 42 synchronization of crossmodal information in heteromodal cortical areas.

43

A number of event-related potential (ERP) studies have examined the temporal aspects
of cross-modal interactions, with the hypothesis that the decrease in reaction time and
facilitation of object recognition should be visible in electrophysiological data. These
studies have found significant activity within several latency windows, with the most

surprising results for audio-visual interactions coming at ~50 ms post-stimulus onset, suggesting extremely early processing of audiovisual interactions (Molholm S *et al.*, 2002). In addition, several ERP studies have also evaluated facilitation of bimodal interactions via an additive model (Besle J et al., 2004). These studies typically have shown amplitude and latency facilitation due to bimodal interactions localized to multimodal cortical areas, as well as suppression of electrophysiological responses with cortical generators in (putatively) unimodal areas.

55

56 A slightly different electrophysiological paradigm for investigating the computational 57 advantages of cross-modal interactions is provided by the steady-state response (SSR), 58 which is the result of entrainment to the physical/spectral properties of a modulated 59 stimulus. This response has been found for both visual and auditory signals and has 60 been used extensively for clinical and diagnostic purposes (Sohmer H et al., 1977). 61 Auditory SSRs are generally elicited by amplitude or frequency modulated signals (e.g. 62 (Luo H et al., 2006)), while visual SSRs are typically elicited by transient high-contrast 63 stimuli such as checkerboard reversals or luminance flicker. Though commonly 64 measured with EEG, the same principles of frequency entrainment to periodic stimuli 65 have been evaluated in MEG as well (Müller MM et al., 1997; Ross B et al., 2000). Ecological stimuli that are temporally extended, and have a quasi-steady-state nature, 66 67 such as speech, can also be modeled via stimuli that approximate the excitation 68 produced by domain-specific information (Grant KW and PF Seitz, 2000). SSRs have a 69 potential further advantage: they can be used to exploit endogenous cortical 70 oscillations. These oscillations are amplified when preferential stimuli (i.e. stimuli that

match the frequency and phase of the endogenous oscillations) constitute the sensory
input (Schroeder CE and P Lakatos, 2009; Schroeder CE *et al.*, 2008; Senkowski D *et al.*, 2008). Oscillatory activity of particular interest occurs in frequency ranges that are
important for relevant behaviors such as speech comprehension, working memory
function and selectional attention (Senkowski D *et al.*, 2008).

76

77 The motivation for the current study was to model an ecologically valid audio-visual 78 interaction, namely speech, using artificial signals that incorporate some critical 79 attributes of a multi-sensory speech. The auditory component of speech consists of the 80 frequency and fine spectral components as well as the envelope – reminiscent of an 81 amplitude-modulated (AM) sinusoidal auditory signal. The speech signal itself shows 82 significant AM activity in the 2 - 16 Hz range (Steeneken HJM and T Houtgast, 1980), 83 and it has been shown that cortical decomposition of the speech envelope is particularly 84 sensitive to frequencies in the range of 4 - 16 Hz. Recent MEG evidence supports this 85 generalization: Luo & Poeppel (Luo H and D Poeppel, 2007) and Howard & Poeppel 86 (Howard MF and D Poeppel, 2010) observed that fluctuations in the speech envelope are associated with intrinsic oscillations in the theta frequency band ($\sim 4 - 8$ Hz). Paired 87 88 with the auditory signal is a visual component in which facial features -- and especially 89 mouth movements -- aid comprehension, especially in noisy environments (Sumby WH 90 and I Pollack, 1954). We crafted stimuli consisting of modulated auditory and visual 91 components within the frequency range of the envelope of speech. By building on 92 results investigating SSRs to auditory and visual stimuli presented alone, we assess the 93 SSR to bimodal audio-visual signals. For this experiment, the visual signal consists of a

94 size-modulated disc (to approximate a mouth opening and closing), and the auditory 95 signal consists of an amplitude-modulated sine wave (to approximate the envelope). We 96 hypothesize that the SSRs elicited by congruent audio-visual signals should be 97 greater than the responses elicited by unimodally modulated auditory or visual stimuli as 98 reflected by the amplitude spectrum at the modulation frequency and the second, third, 99 and fourth harmonics. The increased signal power of the comodulated conditions 100 relative to unimodal conditions might lead to increased activity due to synchrony of 101 different neural populations involved in evaluating the multimodal signal. By 102 manipulating the phase congruence of one modality relative to the other, we additionally 103 aimed to elucidate the online cross-talk between modalities. 104 105 **Materials and Methods** 106 107 Participants: Thirteen right-handed (Oldfield RC, 1971) adult subjects (seven female) 108 with normal hearing and normal or corrected-to-normal vision underwent MEG

109 scanning. One person's data set was excluded from all analyses due to insufficient

110 signal-to-noise ratio for all experimental conditions. Age range was 18-41 (mean 27.08

111 years). Participants were either compensated for their participation or earned course

112 credit in an introductory linguistics course. Presentation of stimuli and biomagnetic

113 recording was performed with the approval of the institutional committee on human

research of the University of Maryland, College Park. Prior to the start of the

115 experiment, written informed consent was obtained from each participant.

117 Stimuli: The experimental stimuli consisted of five types of audio-visual signals 118 presented at two modulation frequencies, for a total of ten signals (Figure 1). The five 119 types were: i) amplitude-modulated sine waves presented concurrently with a static 120 white square on black background; ii) a radius-modulated white disc on black 121 background presented with approximately Gaussian white noise; iii) a radius-modulated 122 disc and an amplitude modulated sine wave at one of three phase relationships (in 123 phase, $\pi/2$ radians out of phase, π radians out of phase). The amplitude-modulated 124 sine waves and radius-modulated discs were modulated at either 2.5 Hz or 3.7 Hz with 125 a modulation depth of 24 percent. These values, a little bit lower than the peak of the 126 modulation spectrum for spoken language, were chosen after extensive piloting 127 revealed that higher visual modulation frequencies were very uncomfortable for 128 participants to view for extended periods of time. Two frequencies were chosen to 129 replicate any effects at different, not harmonically related modulation frequencies. The 130 stimuli were four seconds in duration. For the comodulated conditions, the auditory and 131 visual signal components had the same onset and offset, with the auditory component 132 reaching the maximum value of the modulation envelope first. 133 134 Figure 1 about here

135

Auditory signal components were generated with Matlab (v2007b, The Mathworks,
Natick, MA) and consisted of a sine wave envelope (either 2.5 Hz or 3.7 Hz modulation
frequency) applied to an 800 Hz sine wave carrier signal with 6 ms cos² onset and
offset ramps presented at approximately 65 dB SPL. The signals were sampled at 44.1

140 kHz with 16-bit resolution. Signals were generated using the sine, not the cosine 141 function, to eliminate undesired phase effects on onset responses (see below). Visual 142 signal components were generated using GIMP (www.gimp.org). The radius-modulated 143 white discs were centered on a 640 x 480 pixel black background, and ranged from 2.5° 144 visual angle at the minimum diameter and 4° visual angle for the maximum diameter. 145 The individual frames were compiled into .avi format using VirtualDub 146 (www.virtualdub.org) for presentation. Stimulus timing/frequency was verified with an 147 oscilloscope. The visual components were projected on a screen approximately 30 cm 148 from the participant's nasion. Participants were supine in the MEG scanner for the 149 duration of the experiment.

150

151 Experimental stimuli were presented in nine blocks, with three repetitions per signal per 152 block. Presentation of conditions was randomized within blocks. The experimental 153 materials were passively attended to; no response to the signals was required. In order 154 to maintain vigilance, a distracter task was incorporated into the experiment. An audio-155 visual signal (500 or 1500 ms duration) consisting of a crosshair on a black background 156 combined with approximately Gaussian white noise was used as the target and was 157 pseudorandomly presented with the signals (~17% of total trials). Subjects had to press 158 a button in response to the crosshair/noise target; these trials were excluded from 159 analysis.

160

Delivery: All experimental stimuli were presented using a Dell Optiplex computer with a
 SoundMAX Integrated HD sound card (Analog Devices, Norwood, MA) via Presentation

stimulus presentation software (Neurobehavioral Systems, Inc., Albany, CA). Stimuli
were delivered to the subjects binaurally via Eartone ER3A transducers and nonmagnetic air-tube delivery (Etymotic, Oak Brook, IL). The inter-stimulus interval varied

166 pseudo-randomly between 2500 and 3500 ms.

167

Recording: Data were acquired using a 160-channel whole-head biomagnetometer with
axial gradiometer sensors (KIT System, Kanazawa, Japan). Recording bandwidth was
DC-200 Hz, with a 60 Hz Notch filter, at 1000 Hz sampling rate. Data were noise
reduced using time-shifted PCA (de Cheveigné A and JZ Simon, 2007) trials averaged
offline (artifact rejection ± 2.5 pT), bandpass filtered between .03 - 25 Hz (161 point
Hamming window) and baseline corrected over the 700 ms pre-stimulus interval.

174

175 Data Analysis

The analysis was performed in sensor space, not source space, to stay as close as possible to the recorded data without making source configuration assumptions. All analyses -- pre-experiment localization parameters, waveform assessment, and the calculation of the magnitude and phase of the SSR as well as significance values -were performed in Matlab. Statistical analysis of SSR parameters was evaluated using the statistical and probability distribution functions in Matlab's Statistics Toolbox.

182

Sensor selection from pre-test: Determination of maximally responsive auditory and
 visual channels was performed in separate pre-tests. The auditory pre-test consisted of
 amplitude-modulated sinusoidal signals with 800 Hz sinusoidal carrier signal,

186 modulation frequency 7 Hz, modulation depth 100 percent and 11.3 second duration. 187 The visual pre-test consisted of a checkerboard flicker pattern (Fm = 4 Hz), of 240 188 second duration. The sensor space was divided into quadrants to characterize the 189 auditory response and sextants to characterize the visual response based on the peak 190 and trough field topography expected for each modality as recorded from axial 191 gradiometers (see Figure 1c). Sensor channel designations were anterior temporal 192 (front of head), posterior temporal (rear quadrants/ middle of head) and occipital (back 193 of head overlying occipital lobe). Five channels from source and sink from each sensor 194 division (i.e. ten channels for auditory response and five channels for visual response 195 per hemisphere; 15 channels per hemisphere total) with the maximum measured 196 magnetic field deflection were used for subsequent analyses.

197

198 The auditory pre-test response was characterized using two distinct methods. The first 199 analysis examined the power spectral density (PSD) of the response and selected the 200 channels with the best (strongest) response (Fourier Transform window: 3 to 5 s), at the 201 modulation frequency. The second analysis examined the maximum field deflection of 202 the M100 response (search window: 80 to 130 ms after stimulus onset) and selected 203 the channels with the maximum response amplitude (both source and sink). The pre-204 test visual response was characterized only using the PSD, at twice the modulation 205 frequency (the reversal rate), due to the checkerboard pattern not generating a robust 206 onset response permitting an onset response analysis. Since the data were analyzed in 207 sensor space rather than source space, special care was taken to avoid having 208 posterior temporal and occipital sensors overlap. When posterior temporal and occipital sensors were common to each modality/sensor area, those particular posterior temporal
sensors were replaced by the next non-overlapping posterior temporal sensors.

211

212 Onset response evaluation and PCA: The signal evaluation window (averaged and 213 filtered sensor data) ranged from 700 ms pre-trigger to 3999 ms post-trigger. Onset 214 peak root-mean-square (RMS), RMS latency, magnetic field deflection and magnetic 215 field deflection latency responses corresponding to the M100 (auditory; search window: 216 80 to 130 ms after stimulus onset) and M170 (visual; 145 to 195 ms after stimulus 217 onset) for each hemisphere for each condition were collected and averaged across 218 subjects for each stimulus and were plotted topographically to examine the response. 219 The minimum number of trials averaged was twelve and the maximum number was 220 twenty-seven. Since it is hypothesized that the neurophysiological response primarily 221 reflects processing of the envelope for both signal onset and SSRs, an estimation of the 222 envelope was made using principal components analysis (PCA). The preselected 223 sensors were analyzed using PCA and the envelope estimate was calculated using the 224 absolute value of the Hilbert transform for the first principal component, which explained 225 60 to 80 percent of the total variance, depending on the participant. The channels used 226 for the latency and envelope analysis for the anterior and posterior temporal channels 227 are those from the second analysis (onset analysis described above) of the auditory 228 pre-test data. Congruence between the two sets of pretest channel data was 229 approximately 90%. Occipital channels used were the same as in the pre-test. 230

231 SSR analysis: The magnitude and phase spectra of the SSR were determined using the 232 Fast Fourier Transform (FFT) of the baseline corrected and filtered channel data. The 233 FFT was calculated from stimulus onset (0 ms) to the end of the signal evaluation 234 window (3999 ms). Prior to the calculation of the FT, the data within the signal 235 evaluation window was multiplied by a Kaiser window (length 4000 samples, beta = 13) 236 to minimize the onset and offset responses to the audio-visual signals and minimize 237 spurious frequency contributions. The magnitude of the response was calculated using 238 the RMS of the FT across channels. The phase response was determined by 239 calculating the mean direction as described by Fisher (1996) based on the phase angle 240 of the Fourier transformed data. The across subject response power was determined by 241 calculating the mean of the individual subject power vectors. To determine the across 242 subject phase response, the mean direction of the individual mean directions was 243 calculated. The trials analyzed for the SSR analysis were the same as those analyzed 244 for the onset responses. Figure 2 illustrates waveform recordings for pre-windowed data 245 over the entire analysis frame, the onset responses in detail, and the sensor layout. 246 247 Figure 2 about here

248

SSR cross-modal control analysis: To determine the validity of the sensor selection from the pre-experiment localization, unimodal modulation data were analyzed using the sensors from the other modality. This analysis confirmed that the responses recorded from the unimodal modulation truly reflected that particular modality. This particular analysis does not necessarily indicate that the unimodal visual condition had an effect on the unimodal auditory condition whereas the converse was not true; rather it means
 that the neurophysiological signals generated and recorded were truly present in the
 recorded magnetic field data.

257

258 Across-subject response averaging: The across-subject responses were computed by 259 collecting the individual subject field deflections (source and sink field deflections and 260 RMS) and calculating the mean response amplitudes and the RMS of the subject RMS 261 values. The aggregate waveforms peaks and latencies were characterized in the same 262 search windows as described above. However, the data were not subject to windowing 263 to preserve the onset responses. A similar procedure was used for the Fourier 264 transformed data. Individual subject vectors for response power (squared magnitude) 265 and phase were collected the relevant statistics calculated.

266

267 Statistical analyses: To assess the effect of signal manipulation on the underlying 268 neurophysiological computations, the mean latency and peak values (for both magnetic 269 field sink and source deflections and RMS) for onset responses were analyzed using 270 mixed measures ANOVA (SPSS 16.0, SPSS Inc., Chicago, IL). A full factorial design 271 was employed, with amplitude (peak RMS and maximum response) and latency as the 272 dependent measures and Hemisphere (RH vs. LH), Frequency (2.5 Hz vs. 3.7 Hz) and 273 Phase (in phase, $\pi/2$ radians out of phase, π radians out of phase) as factors. Planned 274 comparisons using Wilcoxon sign rank tests compared unimodal modulation against 275 each comodulated condition for i) peak RMS, ii) peak RMS latency, iii) peak source and 276 sink deflection, iv) source and sink deflection latency, v) envelope onset period, vi) peak of envelope onset and vi) envelope periodicity if and only if the ANOVA results werefound to be significant.

279

280 The significance of the SSR amplitude at a specific frequency was analyzed by 281 performing an F test on the squared RMS (power) of the Fourier transformed data 282 (Dobie & Wilson 1996). The signal evaluation window used gave a frequency resolution 283 of 0.25 Hz and gave the exact response at Fm = 2.5 Hz, but not at 3.7 Hz. To evaluate 284 the response at Fm = 3.7 Hz, the bin next closest in frequency (3.75 Hz) was used. 285 The significance of the phase of the response was assessed using Rayleigh's phase 286 coherence test on the mean direction (Fisher NI, 1996). Individual subject responses at 287 each modulation frequency for each condition were assessed using an F test to 288 determine if the response was significant and whether or not a particular subject should 289 be excluded due to lack of a response or exhibiting a response other than at the 290 modulation frequencies and harmonics of interest. For the across-subject data, F tests 291 were performed on the power of the SSR at the modulation frequency, two 292 subharmonics in the delta band, one harmonic in the theta band and one in the alpha 293 band (see e.g., Jones & Powell 1970; Senkowski et al. 2008 for review of frequency 294 band descriptions). The power at individual harmonic components of the modulation 295 frequency in different frequency bands across conditions was compared using Wilcoxon 296 sign rank tests (Matlab v7). Two sets of sign rank tests were performed: the first 297 compared the mean unimodal modulation magnitudes against the mean comodal 298 modulation magnitudes for a given sensor area (e.g. LH anterior temporal unimodal 299 auditory vs. LH anterior temporal comodal, phi = π) and the second compared the

300 comodulated conditions (e.g. RH occipital, phi = zero vs. RH occipital, phi = $\pi/2$. A 301 mixed effects ANOVA implemented in R (Baayen 2008; R Development Core Team 302 2008) assessed any possible differences in modulation frequency and hemisphere. 303

304 Dipole estimation: SSR source estimation was performed on data from seven subjects 305 who exhibited a SSR response. Since we did not have structural MR images for our 306 participants, source estimations were not anatomically constrained. Single equivalent 307 current dipole estimates with a GOF < 80% and not localized to the hemisphere from 308 which the channels were selected were excluded from subsequent statistical analyses. 309 A simple spherical head model was used to determine the source of the SSR (x,y,z)310 axes) as well as the dipole angles (theta and phi) using 35 sensors per hemisphere with 311 the greatest PSD (ten channels each from anterior and temporal sensor divisions; 312 fifteen from occipital sensor divisions per hemisphere). The sensors selected came 313 from both the auditory and visual sensor division described previously. To perform the fit 314 to the desired components found significant by the F test (see Results), the real part of 315 the Fourier transformed data for the component of interest was multiplied by 316 $\cos(2^{*}\pi^{*}F^{*}t)$ and the imaginary part of the data by $\sin(2^{*}\pi^{*}F^{*}t)$, where F and t are the 317 frequency of interest and the time vector, respectively. The resulting vectors were then 318 algebraically added and fits were performed on the peaks of the variance of the 319 resulting sinusoidal waveform. Peaks corresponding to the magnetic field topography 320 reflected by the response magnitudes were used for the source estimation (see 321 *Results*). Dipole estimations from a single peak was recorded and entered for 322 subsequent statistical analysis. Statistical significance of the values of x,y,z, theta and

phi were assessed using a mixed effects ANOVA in R using the 'languageR' statistical
package. Wilcoxon tests were performed on the values of theta and phi both across
and within subjects to assess any potential differences in the source orientation of the
SSR (R 2.81).

327

328 Results

329

330 Figure 3 illustrates the response for one participant for both unimodal modulation 331 conditions, Fm = 2.5 Hz. This characteristic pattern demonstrates that the majority of 332 evoked activity occurs primarily at the modulation frequency. The least amount of 333 response power was observed to be in the sensors overlying the anterior temporal lobe 334 (analyzed for auditory alone condition only), with greater magnitude in the posterior 335 temporal sensors for auditory alone and for the occipital sensors in visual alone 336 unimodal conditions. The overall response fit well with that of a prototypical SSR, 337 namely the response was elicited robustly at the modulation frequency and occasionally 338 visible at the first few harmonics. 339 340 Figure 3 about here 341

The response topography for a representative subject is illustrated in Figure 4. The complex-valued magnetic field response profile reflects the SSR response power at the modulation frequency for each condition, Fm = 2.5 Hz, as measured at one peak in the sinusoidal waveform used for dipole localization (discussion below). The topography 346 reflects the whole-head response power measured for each condition. For the unimodal 347 modulation conditions, the whole-head response resembles the topography of a visual 348 response. This is congruent with the observation that the visual response was greater 349 than the auditory response in the unimodal modulation conditions (see Figure 3 and 350 Figure 6). The bimodal condition magnetic field topographies reflect the combined 351 auditory and visual cortical computations underlying their generation. The increasing 352 contribution of auditory cortex to cortical processing of the bimodal signals can be seen 353 in the spatial configuration of the magnetic fields recorded. The contribution of the 354 auditory cortex can be observed most clearly when the signal component envelopes are 355 initially orthogonal to one another (Figure 4d; $\Phi = \pi/2$). Magnetic field sink-source 356 orientation reverses at each peak for this condition. 357 358 Figure 4 about here 359

360 Across-Subject Power Analysis

361 Figure 5 displays the across subject response power for Fm = 3.7 Hz, plotted with a 362 linear scale for frequency and a logarithmic scale for response power, shown here for 363 right hemisphere sensors only. Across conditions it is evident that the anterior channels 364 capture the underlying SSR activity in a less compelling manner; posterior temporal and 365 occipital channels, on the other hand, reveal very clear patterns. Response power was 366 concentrated at the modulation frequency and the second harmonic, with some activity 367 centered also around 10 Hz, as was clear for the representative subject shown above. 368 Figure 5 shows the grand-averaged response power across subjects for the right

369	hemisphere only, Φ = 0, Φ = $\pi/2$ and Φ = $\pi,$ Fm = 3.7 Hz. Response power significance
370	for all bimodal conditions (as determined by Wilcoxon sign-rank tests) compared to the
371	unimodal modulation conditions show that the responses are significantly greater in
372	bimodal than unimodal responses at the frequencies found significant by the ANOVA.
373	
374	Figure 5 about here
375	
376	Several observations merit highlighting. First, the majority of the activity is reflected in
377	the sensors overlying the posterior temporal lobes and occipital lobes. Second, for the
378	AV comodulated condition in which the signal envelopes are at the same initial phase,
379	the response power is greatest at the modulation frequency, localized to the sensors
380	overlying the posterior temporal lobes. Third, as the difference in the relative phase
381	increases, the response power decreases, although the response is still greater than
382	that of unimodal modulation condition (see Figure 6). This change is indexed primarily
383	by the response power as measured at the modulation frequency.
384	
385	Statistical summary
386	The significance of the SSR was calculated at the modulation frequency, as well as two
387	subharmonics, and the second and third harmonics. Significance was determined by
388	means of an <i>F</i> test on the power of the SSR at each particular frequency as described
389	by Dobie and Wilson ((Dobie RA and MJ Wilson, 1996) - see Methods). All subjects
390	elicited a statistically significant response for the SSR at each envelope modulation
391	frequency. Within-subject response significance was restricted to evaluation at the

392	modulation frequency (see <i>Methods</i>) with degrees of freedom (df) df1 = 2, df2 = 12 and					
393	α = 0.05. The across-subject significance for subharmonics was assessed using <i>df</i> =					
394	2,4 and significance for the modulation frequency and second and third harmonics were					
395	assessed using $df = 2,12$.					
396						
397	SSR power at subharmonics was not found to be statistically significant. Statistically					
398	significant responses were observed at the modulation frequency, as well as second					
399	and third harmonics. The difference between the observed statistical significance for					
400	subharmonics and the second and third harmonics may be attributable to the decreased					
401	degrees of freedom for df2.					
402						
403	Results of Rayleigh's test on the mean direction of the SSR vectors (at the frequencies					
404	observed to be significant by the <i>F</i> test) found the phase angle directions to be					
405	statistically significant at α = 0.05					
406						
407	Figure 6 about here					
408						
409	SSR power comparisons					
410	Table 1 summarizes the SSR power for each modulation frequency, modulation					
411	condition and sensor/cortical area as well as the interactions found to be significantly					
412	significant as a result of Wilcoxon sign-rank tests (see Methods). For both envelope					
413	modulation frequencies, several statistically significant responses are held in common.					
414	First, both modulation frequencies exhibit statistically significant responses power at the					

415 modulation frequency for all comodulated conditions and this interaction is largely 416 limited to the sensors overlying the posterior temporal lobe for both hemispheres. 417 Second, there were significant interactions at the second harmonic for $\Phi = 0$ and $\Phi = \pi$: 418 both modulation frequencies held this interaction in common in the LH sensors overlying 419 the posterior temporal lobe. One last interaction was common to both modulation 420 frequencies for the third harmonic for $\Phi = 0$ in the LH sensors overlying the occipital 421 lobe. Several other statistically significant interactions were found to be unique to each 422 modulation frequency; these perhaps inconsistent interactions may be a result of data 423 variance (see *Discussion*). No statistical difference was observed for SSR power 424 between the three bimodal conditions. Linear mixed effects models with modulation 425 frequency and hemisphere as factors found no significant statistical interactions. 426 Wilcoxon sign-rank tests were performed on the incidental power centered around 10 427 Hz to determine if it was significant. Results of the tests across conditions yielded no 428 significant results. 429 430 Table 1 about here 431 432 Dipole source estimation 433 Based on the literature, we assumed several cortical locations to be implicated for each 434 condition of the experiment. The response for the unimodal auditory condition should be 435 localized to early auditory cortex (Ross B et al., 2000) the response for the unimodal

436 visual condition to occipital visual cortex (Müller MM et al., 1997) For the bimodal

437 conditions, we assumed an SSR source localization to superior temporal sulcus or

438 superior parietal lobule, both regions previously suggested to underlie AV integration 439 (e.g. (Molholm S et al., 2006)). However, equivalent current dipole estimation was 440 largely unsuccessful. The majority of localization estimates had a goodness of fit values 441 < 80% and for estimates that met our criteria, localization was more successful for the 442 LH than for the RH. When dipole localization was successful at all, dipoles were 443 localized to the LH, corresponding perhaps to the parietal lobe area (Molholm S et al., 444 2006). Because this study is explicitly about the detection and modulation of a specific 445 electrophysiological response, the SSR, we did not design the study with source 446 localization in mind (and indeed, we do not have subject structural MRIs to permit 447 adequate source localization). Although the unimodal SSR responses are well 448 characterized in the literature and their superior temporal and occipital localizations not 449 controversial, it would be helpful to be able to localize the bimodal SSR here, but we 450 were not able to with sufficient accuracy.

451

452 PCA (onset) data for envelopes and response amplitudes

453 Statistical evaluation of the first principal component yielded no significant interactions. 454 However, several potential interesting possibilities warranting further studies were 455 observed. First, the rise time of the onset response may vary with the relative phase of 456 the components. When both signal component envelopes are completely synchronized, 457 the rise time appeared to be fastest; when completely desynchronized the rise time 458 appeared to be slowest. There are potential differences between hemispheres and 459 sensor/cortical divisions. Additionally, the initial part of the onset response seems to 460 reflect the nature of the signals presented: there is initial sinusoidal activity that reflects

the both the AM envelope of the auditory signal component as well as the increases and
decreases in the radius of the visual signal component. Occurring prior to signal
entrainment and the SSR, this onset activity has roughly the same duration as the
period of the modulation frequency.

465

466 As with the PCA, no statistically significant interactions were observed for the mean 467 magnetic field deflections. However, as for PCA, there were several interactions that warrant follow-up. Though no latency or peak field values were found to be significant. 468 469 the spatial configuration of the observed magnetic fields may vary based the type of 470 modulation (unimodal vs. bimodal) as well as the phase of the signal components. 471 Contributions from the posterior temporal lobe/sensors result in the magnetic field 472 displaying a more 'auditory' or heteromodal character. As with the previously described 473 SSR topographies, the observed magnetic field was a mixture of auditory and visual 474 magnetic field topographies.

475

476 Discussion

477

The study examines the steady state response properties to multisensory signals. We demonstrate that the oscillatory activity in the human auditory and visual cortices entrains to periodic stimuli as reflected in increased power in the SSR at the modulation frequencies. Given that auditory SSRs and visual SSRs have been utilized as robust diagnostic tools, this result is not surprising. However, the types SSR studies differ from the current experiment in three important ways. First, our visual stimulus differed from the usual transient (checkerboard or flicker) stimuli used in SSVEPs. We used gradually
growing and shrinking discs in an effort to approximate the opening and closing of a
human mouth producing speech and were successful in eliciting steady state responses
at the modulation frequencies employed. The modulation frequencies chosen here
were based on previous findings regarding the temporal modulation of the speech
envelope (Chandrasekaran C et al., 2009)

490

491 A second important difference is the use of a *combination* of auditory and visual signals 492 to elicit the SSR and to evaluate whether entrainment to comodulated, multisensory 493 stimuli differs from unimodally modulated stimuli. We find that the multisensory, 494 comodulated signals elicit a SSR with greater power than either unimodal signal. This 495 experiment was designed to model the potential tracking of speech amplitude 496 envelopes at an approximately syllabic rate (~ 4 Hz). The neurophysiological 497 mechanisms and behavioral consequences of crossmodal integration may be related to 498 that of the auditory phenomenon known as comodulation masking release (CMR) 499 (Grant KW and PF Seitz, 2000). This additional power at the modulation frequency may 500 reflect the neural advantage behind the perceptual boost gained from the release from 501 masking in bimodal speech.

502

503 In CMR, an auditory signal can be detected even at poor signal to noise ratios due to 504 the presence of a comodulated stimulus. The degree of release from masking is 505 greatest when the masker bandwidth is large and has a high spectrum level, the 506 modulation frequency of the signal is low, has a high modulation depth and regular envelope (Verhey JL et al., 2003). This has obvious parallels with audio-visual speech
signals where background noise (approximately white noise) can mask the auditory
signal of interest, the speech envelope is modulated at low frequencies and has a
relatively high modulation depth. The mechanisms involved for such an auditory-only
process could logically be extended not only to other modalities, but also to crossmodal
computations.

513

514 Third, this study investigates levels of asynchrony to evaluate the *tolerance* of this SSR 515 to shifts in phase between the modalities. We did not, however, show statistically 516 significant differences in the power of the responses at the modulation frequency as a 517 function of the phase shift between the signals in the bimodal conditions.

518

519 We hypothesize that the initial differences in phase give rise to distinct representations 520 of the bimodal signals. For the condition in which both signal component envelopes are 521 completely synchronous, the AV signal is computed (integrated) as being a single 522 object. When the starting phases are orthogonal to each other, the AV signal is 523 alternately evaluated as being either two signals or one signal, as the initial phase 524 difference causes the component envelopes to synchronize and desynchronize over the 525 duration of the signal. And when the signal component envelopes are π radians out of 526 phase, each component is evaluated as being a separate object. When the signal 527 component envelopes are completely synchronized, this is analogous to a highly 528 correlated statistical regularity in the bimodal signal. As the signal components are 529 desynchronized, these correlations and redundancies in the bimodal signal decrease,

530 modifying the processing and representation of the percept. Chandrasekaran et al.

(2009) employed bimodal speech stimuli (with no phase incongruities) and observed (i)
a temporal correspondence between mouth opening and the auditory signal component
envelope and (ii) mouth openings and vocal envelopes are modulated in the 2 – 7 Hz
frequency range.

535

536 While we have demonstrated that the feasibility of eliciting a bimodal SSR and that the 537 SSR indexes congruency, the indexing performed by the SSR, at least with the 538 conditions in this particular experiment, are limited. For all bimodal conditions we 539 observed that their response power was greater than that of the unimodal conditions 540 and that the further the signal components were separated in phase, the response 541 power decreased. However the congruency indexed by the phase separation may have 542 practical limits. There is evidence that integration of bimodal signals, with the auditory 543 signal component leading, takes place within a 40-60 ms duration window (van 544 Wassenhove V et al., 2007). For the modulation frequencies employed in this 545 experiment, the incongruity between signal components did not fall within this 546 integration window. It is entirely possible that the incongruity tolerance is dependent on 547 the modulation frequency. Higher envelope modulation rates (e.g. 7-11 Hz) will yield 548 phase separations that can test the tolerance between signal components. A related 549 issue is to sample more phase separation values around the entire unit circle. One 550 possible hypothesis is that the representation of the phase separation will be symmetric 551 (except when both signal envelopes are completely synchronized), i.e. the response 552 power for a phase separation of $\pi/2$ radians and $3\pi/2$ radians will be represented

553 equally. The indexing of signal component congruity might also be dependent on which 554 component reaches the maximum of the envelope first. It has been shown that when 555 visual information precedes auditory information, signal detection and comprehension 556 increases (Senkowski D et al., 2008; van Wassenhove V et al., 2007). In the current 557 study, for the asynchronous bimodal conditions, the auditory component of the signal 558 reached the maximum of the modulation envelope first. Future studies could investigate 559 the converse situation, where the modulation envelope of the visual signal component 560 reaches the maximum prior to that of the auditory component. Further iterations of this 561 experimental paradigm could investigate a combination of these factors: the impact 562 modulation frequency and phase has on bimodal integration, using a wider range of 563 phase separations and determining to what extent which signal component reaches the 564 maximum envelope value first affects indexing of congruity.

565

566 Subsequent experiments can improve on the paradigm introduced here in several 567 important ways. First, the modulation depth of the auditory signal component might be 568 made more variable. In this study, we aimed to have as much congruency as possible 569 between the auditory and visual components of the signal. To create a more 570 ecologically valid AV signal, the modulation depth should correspond to the conditions 571 occurring in natural human speech, where the mouth opens and closes fully (modulation 572 depth ranging from 0 to 100 per cent). Secondly, the analysis of the signals could be 573 improved in several ways. The amount of trials analyzed were sufficient to characterize 574 the SSR, but not the signal onsets. It may be possible to analyze both the SSR and the 575 signal onsets using PCA by increasing the number of trials while decreasing the trial

576 duration. The SSR may still be recovered by concatenating the trials (Xiang 2008) while 577 to recover the onset information, the trials can be averaged, as with the analysis of 578 M100 or M170 responses. Several potential hypotheses suggest that conducting a dual 579 analysis in this manner would be potentially useful. First, it has been previously argued 580 that the RH shows a slight preference in extracting envelope information in auditory 581 processing (Luo H and D Poeppel, 2007) even though we found no difference in power 582 between hemispheres. The LH might play a larger role in analyzing the onset 583 responses and the increase in trial numbers might provide a more clear set of data to 584 analyze.

585

586 An additional way the signal could be improved is by using ellipsoids rather than circles 587 to simulate mouth movements. This shape is not only much closer to that of a human 588 mouth, but by modulating only the shorter of the radii, this would yield a more natural 589 modulation motion. There was also a technical issue in implementing the current study. 590 We desired to use modulation rates in the 4 - 16 Hz range; however these modulation 591 rates caused shape and color change effects in the visual signal component. Pilot 592 versions of the current study found that these rates and the perceptual effects they 593 caused made the experiment undoable, except at slower modulation rates. Ultimately, 594 we chose modulation rates that agreed with previous research and that did not induce 595 undesired perceptual effects. The human visual system can track changes faster than 596 those of 16 Hz without such issues occurring and thus there may be a local vs. global 597 problem when using circular stimuli as employed in this study. Using the elliptical shape 598 may reduce these effects and allow for the creation and investigation of more

ecologically valid signals. Adding 'jitter' or noise to the signal component envelopes
may also yield a yet more ecologically valid set of stimuli for further experimentation.
This adds the variability inherent in ecological speech, while retaining the modulation
information of the signal component envelopes.

603

In summary, we demonstrate – to our knowledge for the first time -- that an
experimental technique previously solely applied to unimodal signals, the SSR, can be
applied to signals of a bimodal nature. Furthermore, the paradigm reported (as well as
potential modifications) yields a useful index of signal component congruity that can be
applied to the study of speech and other ecologically valid crossmodal interactions.

009

610 Acknowledgements

611

612 This project originated with a series of important discussions with Ken W. Grant 613 (Auditory-Visual Speech Recognition Laboratory, Army Audiology and Speech Center, 614 Walter Reed Army Medical Center). The authors would like to thank him for his 615 extensive contributions to the conception of this work. The authors would like to thank 616 Mary F. Howard and Philip J. Monahan for critical reviews of this manuscript. We would 617 also like to thank Jonathan Z. Simon for outstanding analytical assistance of the 618 experimental data and Jeff Walker for technical assistance in data collection. This work 619 was supported by the National Institutes of Health (2R01DC05660 to DP) and the 620 National Institute on Deafness and Other Communication Disorders of the National 621 Institutes of Health (Training Grant DC-00046 support to JJIII and AER).

623

624 Amedi A, von Kriegstein K, van Atteveldt NM, Beauchamp MS, Naumer MJ (2005) 625 Functional imaging of human crossmodal identification and object recognition. 626 Exp Brain Res 166: 559-571. 627 Baumann O, Greenlee MW (2007) Neural Correates of Coherent Audiovisual Motion 628 Perception. Cereb Cortex 17: 1433-1443. 629 Besle J, Fort A, Delpuech C, Giard MH (2004) Bimodal speech: early suppressive visual 630 effects in human auditory cortex. Eur J Neurosci 20: 2225-2234. 631 Calvert GA, Brammer MJ, Bullmore ET, Campbell R, Iversen SD, David AS (1999) 632 Response amplification in sensory-specific cortices during crossmodal binding. 633 Neuroreport 10: 2619-2623. 634 Calvert GA, Campbell R, Brammer MJ (2000) Evidence from functional magnetic 635 resonance imaging of crossmodal binding in the human heteromodal cortex. 636 Current Biology 10: 649-657. 637 Calvert GA, Hansen PC, Iversen SD, Brammer MJ (2001) Detection of Audio-Visual 638 Integration Sites in Humans by Application of Electrophysiological Criteria to the 639 BOLD effect. NeuroImage 14: 427-438. 640 Campbell R (2008) The processing of audio-visual speech: empirical and neural bases. 641 Philos Trans R Soc Lond B Biol Sci 363: 1001-1010. 642 Chandrasekaran C, Trubanova A, Stillittano S, Caplier A, Ghazanfar AA (2009) The 643 Natural Statistics of Audiovisual Speech. PLoS Comput Biol 5: e1000436. 644 de Cheveigné A, Simon JZ (2007) Denoising based on time-shift PCA. J Neurosci 645 Methods 165: 297-305. 646 Dobie RA, Wilson MJ (1996) A comparison of t test, F test, and coherence methods of 647 detecting steady-state auditory-evoked potentials, distortion product otoacoustic 648 emissions, or other sinusoids. J Acoust Soc Am 100: 2236-2246. 649 Driver J, Spence C (1998) Crossmodal attention. Curr Opin Neurobiol 8: 245-253. 650 Fisher NI (1996) Statistical Analysis of Circular Data. Cambridge: Cambridge University 651 Press. 652 Fort A, Delpuech C, Pernier J, Giard M-H (2002) Dynamics of Cortico-subcortical Cross-653 modal Operations Involved in Audio-visual Object Detection in Humans. Cereb 654 Cortex 12: 1031-1039. 655 Grant KW. Seitz PF (2000) The use of visible speech cues for improving auditory 656 detection of spoken sentences. J Acoust Soc Am 108: 1197-1208. 657 Hershenson M (1962) Reaction time as a measure of intersensory facilitation. J Exp 658 Psychol 63: 289. Howard MF, Poeppel D (2010) Discrimination of Speech Stimuli Based on Neuronal 659 660 Response Phase Patterns Depends On Acoustics But Not Comprehension. J 661 Neurophysiol. 662 Jones EG, Powell TP (1970) An anatomical study of converging sensory pathways 663 within the cerebral cortex of the monkey. Brain 93: 793-820. 664 Kayser C, Petkov CI, Logothetis N, K. (2008) Visual Modulation of Neurons in Auditory 665 Cortex. Cereb Cortex 18: 1560-1574.

- Kelly SP, Gomez-Ramirez M, Foxe JJ (2008) Spatial Attention Modulates Initial Afferent
 Activity in Human Primary Visual Cortex. Cereb Cortex 18: 2629-2636.
- Lakatos P, Karmos G, Mehta AD, Ulbert I, Schroeder CE (2008) Entrainment of
 Neuronal Oscillations as a Mechanism of Attentional Selection. Science 320:
 110.
- Lalor EC, Kelly SP, Pearlmutter BA, Reilly RB, Foxe JJ (2007) Isolating endogenous
 visuo-spatial attentional effects using the novel visual-evoked spread spectrum
 analysis (VESPA) technique. Eur J Neurosci 26: 3536--3542.
- Luo H, Poeppel D (2007) Phase Patterns of Neuronal Responses Reliably Discriminate
 Speech in Human Auditory Cortex. Neuron 54: 1001-1010.
- Luo H, Wang Y, Poeppel D, Simon JZ (2006) Concurrent Encoding of Frequency and
 Amplitude Modulation in Human Auditory Cortex: MEG Evidence. J Neurophysiol
 96: 2712-2723.
- 679 Macaluso E, Driver J (2005) Multisensory spatial interactions: a window onto functional 680 integration in the human brain. Trends Neurosci 28: 264-271.
- 681 Mesulam MM (1998) From sensation to cognition. Brain 121: 1013-1052.
- 682 Miller BT, D'Esposito M (2005) Searching for "the Top" in Top-Down Control. Neuron 683 48: 535-538.
- Molholm S, Martinez A, Shpaner M, Foxe JJ (2007) Object-based attention is
 multisensory: co-activation of an object's representations in ignored sensory
 modalities. Eur J Neurosci 26: 499-509.
- Molholm S, Ritter W, Javitt DC, Foxe JJ (2004) Multisensory Visual-Auditory Object
 Recognition in Humans: a High-density Electrical Mapping Study. Cereb Cortex
 14: 452-465.
- Molholm S, Ritter W, Murray MM, Javitt DC, Schroeder CE, Foxe JJ (2002)
 Multisensory auditory-visual interactions during early sensory processing in
 humans: a high-density electrical mapping study. Cognitive Brain Research 14:
 115-128.
- Molholm S, Sehatpour P, Mehta AD, Shpaner M, Gomez-Ramirez M, Ortigue S, Dyke
 JP, Schwartz TH, Foxe JJ (2006) Audio-Visual Multisensory Integration in
 Superiour Parietal Lobule Revealed by Human Intracranial Recordings. J
 Neurophysiol 96: 721-729.
- 698 Müller MM, Teder W, Hillyard SA (1997) Magnetoencephalographic recording of 699 steadystate visual evoked cortical activity. Brain Topography 9: 163-168.
- Murray MM, Foxe JJ, Wylie GR (2005) The brain uses single-trial multisensory
 memories to discriminate without awareness. NeuroImage 27: 473-478.
- Oldfield RC (1971) The assessment and analysis of handedness: the Edinburgh
 inventory. Neuropsychologia 9: 97-113.
- Olson IR, Gatenby JC, Gore JC (2002) A comparison of bound and unbound audio visual information processing in the human cerebral cortex. Cognitive Brain
 Research 14: 129-138.
- Ross B, Borgmann C, Draganova R, Roberts LE, Pantev C (2000) A high-precision
 magnetoencephalographic study of human auditory steady-state responses to
 amplitude modulated tones. J Acoust Soc Am 108: 679-691.
- Schroeder CE, Foxe J (2005) Multisensory contributions to low-level, 'unisensory'
 processing. Curr Opin Neurobiol 15: 454-458.

- Schroeder CE, Lakatos P (2009) Low-frequency neuronal oscillations as instruments of
 sensory selection. Trends Neurosci 32: 9-18.
- Schroeder CE, Lakatos P, Kajikawa Y, Partan S, Puce A (2008) Neuronal oscillaions
 and visual amplification of speech. Trends Cogn Sci 12: 106-113.
- 716Senkowski D, Molholm S, Gomez-Ramirez M, Foxe JJ (2006) Oscillatory Beta Activity717Predicts Response Speed during a Multisensory Audiovisual Reaction Time
- 718Task: A High-Density Electrical Mapping Study. Cereb Cortex 16: 1556-1565.
- Senkowski D, Schneider TR, Foxe JJ, Engel AK (2008) Crossmodal binding through
 neural coherence: implications for multisensory processing. Trends Neurosci 31:
 401-409.
- Sohmer H, Pratt H, Kinarti R (1977) Sources of frequency following response (FFR) in
 man. Electroencephalogr Clin Neurophsyiol 42: 656-664.
- Steeneken HJM, Houtgast T (1980) A physical method for measuring speech transmission quality. J Acoust Soc Am 67: 318-326.
- Stein BE, Meredith MA, Huneycutt WS, McDade L (1989) Behavioral Indices of
 Multisensory Integration: Orientation to Visual Cues is Affected by Auditory
 Stimuli. J Cogn Neurosci 1: 12-24.
- Sumby WH, Pollack I (1954) Visual Contribution to Speech Intelligibility in Noise. J
 Acoust Soc Am 26: 212-215.
- van Wassenhove V, Grant KW, Poeppel D (2007) Temporal window of integration in auditory-visual speech perception. Neuropsychologia 45: 598-607.
- Verhey JL, Pressnitzer D, Winter IM (2003) The psychophysics and physiology of
 comodulation masking release. Exp Brain Res 153: 405-417.
- 735
- 736

738 Figure captions

740	Figure 1a. Schematic of stimuli employed in experiment. Upper panel illustrates the						
741	movement of the visual signal component throughout the duration of stimulus (4						
742	seconds – see Methods for details). Lower panel illustrates the auditory signal						
743	component for the duration of the stimulus. The stimuli were presented at one of two						
744	modulation frequencies (Fm = 2.5 and 3.7 Hz), modulation depth was 24 per cent,						
745	carrier frequency for the auditory signal component was 800 Hz. Synchronous						
746	condition is pictured: maximum radius of circle corresponds to maximum envelope value						
747	for auditory component.						
748							
749	Figure 1b. Visual signal components were centered on a 640 x 480 black background						
750	and presented from 2.5° - 4° visual angle approximately 30 cm from the subjects'						
751	nasion. As in Fig 1a, maximum radius of the circle corresponds to maximum envelope						
752	value for auditory component when both component envelopes are synchronized.						
753							
754	Figure 1c. Division of magnetoencephalographic sensors. Top panel shows division of						
755	auditory sensors for experimental pre-test; bottom panel shows sensor division for						
756	visual pre-test. Sensor division was based on expected field topography for auditory						
757	and visual cortical responses recorded from axial gradiometer sensors (see Methods for						
758	details). Sensor designation is as follows: A = anterior temporal sensors, P = posterior						
759	temporal sensors, O = occipital sensors. Placement of letters roughly corresponds to						
760	the locations of the sensors selected for the analysis of the experimental data.						

762	Figure 2a. Butterfly plot of MEG waveform pre-windowing from a single subject (see						
763	Methods for details). Recorded magnetic field deflections are in black, root-mean-						
764	square (RMS) in red. This illustrates the recorded field deflections (onset and SSR)						
765	prior to multiplication by the Kaiser window in the post-signal onset time domain.						
766							
767	Figure 2b. Magnification of data shown in panel 2a, focusing on the recorded onset						
768	response (prior to multiplication by Kaiser window in the time domain). This illustrates						
769	the necessity of minimizing the onset responses to avoid contamination of the						
770	evaluation of the steady-state portion of the signal.						
771							
772	Figure 2c. Sensor configuration of whole-head biomagnetometer. The waveforms						
773	displayed are from the steady-state portion of the response from 2810 to 3010 ms post-						
774	stimulus onset.						
775							
776	Figure 3. Single subject response power plots. Shown is the response at $Fm = 2.5 Hz$,						
777	for each of the unimodal modulation conditions Left column plots the recorded response						
778	for the LH, right column the response for the RH. Top row illustrates the response from						
779	the anterior temporal sensors, middle row posterior temporal sensors (these rows						
780	correspond to unimodal A condition), bottom row occipital sensors (corresponds to						
781	unimodal V condition). Gray shading highlights the response power at modulation						
782	frequency and second and third harmonics. Responses were found to be significant for						

Figure 4a. Complex-valued magnetic field topography for a single subject, Fm = 2.5 Hz, unimodal auditory modulation. The measured response topography resembles that of a visual response, possibly because the occipital response power is greater than that of the temporal response overall (for comparison, see Figure 3).

789

Figure 4b. Complex-valued magnetic field topography, unimodal visual modulation. Asin (a), the observed power mirrors that of a primarily visual response.

792

Figure 4c – Figure 4e. Complex valued magnetic field topography for the bimodal conditions. The sink/source pattern shows increasing auditory cortex contribution to bimodal processing observed in the magnetic field spatial distribution. Panel (c) plots the response for $\Phi = 0$, panel (d), $\Phi = \pi/2$, panel (e) $\Phi = \pi$.

797

798 Figure 5. Across subject response power for all three bimodal conditions, Fm = 3.7 Hz. 799 Mean response power across subjects for RH sensors only. Frequency (Hz) is plotted on a linear scale; response power (fT^2/Hz) is plotted on a logarithmic scale. Grav 800 801 shading indicates power at the modulation frequency and second harmonic. SSR power 802 peaked at the modulation frequency and the second harmonic, with some incident 803 activity centered around 10 Hz. Response power is concentrated in the sensors 804 overlying the posterior temporal and occipital lobes (middle and bottom rows). 805 Logarithmic plotting on the y-axis is employed to give a sense of the range of the

response power for all sensor areas; however, it skews the data somewhat in that the
powers found to be statistically significant by the *F* test do not appear to be significant.

- 809 Figure 6. Magnitude response plots at modulation frequency for each modulation
- 810 frequency (top row shows Fm = 2.5 Hz, bottom row shows Fm = 3.7 Hz) and each
- 811 hemisphere (left panels show LH, right panels show RH), grouped by experimental
- 812 condition. Magnitude is displayed logarithmically on the y-axis (fT²/Hz), and gray
- 813 shading indicates sensor division. Response power for comodulated conditions is
- 814 greater than response power for unimodal conditions. Greatest response power is seen
- 815 in the posterior temporal sensors when the envelopes are initially orthogonal ($\Phi = \pi/2$).







C auditory sensor division



visual sensor division

b

















Fm = 2.5 Hz									
Con	Unimodal		Φ = 0		Φ = π/2		Φ = π		
Hem	LH	RH	LH	RH	LH	RH	LH	RH	
Ant Tem	Fm=3.253e1 5 2nd=9.985e1 4 3rd=1.783e1 5	Fm=2.750e1 5 2nd=1.180e1 5 3rd=4.888e1 4	Fm=1.136e16 2nd=1.294e1 5 3rd=8.316e14	Fm=2.866e16 2nd=2.797e1 5 3rd=1.092e15	Fm=1.077e1 6 2nd=6.617e1 4 3rd=8.486e1 4	<i>Fm=4.741e1</i> 6 2nd=1.305e1 5 3rd=8.950e1 4	Fm=1.043e16 2nd=8.746e1 4 3rd=1.409e15	<i>Fm=1.494e16</i> 2nd=1.827e15 3rd=3.193e15	
Pos Tem	Fm=3.630e1 5 2nd=1.539e1 5 3rd=1.199e1 6	Fm=2.750e1 5 2nd=1.577e1 5 3rd=5.518e1 6	Fm=2.739e1 6 2nd=5.645e1 5 3rd=7.867e1 6	Fm=5.489e1 6 2nd=1.298e1 6 3rd=2.005e16	<i>Fm=8.04e16</i> 2nd=1.676e1 5 3rd=6.539e1 6	<i>Fm=5.112e1</i> 6 2nd=1.581e1 5 3rd=2.867e1 6	Fm=4.281e1 6 2nd=4.668e1 5 3rd=3.021e16	Fm=2.149e16 2nd=4.7875e1 5 3rd=4.981e16	
Occ	Fm=3.654e1 6 2nd=4.213e1 5 3rd=1.910e1 6	Fm=2.760e1 6 2nd=4.772e1 5 3rd=2.160e1 6	Fm=2.821e16 2nd=8.426e1 5 <i>3rd=1.068e1</i> 7	Fm=4.403e16 2nd=6.545e1 5 3rd=3.176e16	Fm=2.376e1 6 2nd=3.606e1 5 3rd=1.291e1 7	Fm=5.317e1 6 2nd=9.348e1 5 3rd=4.404e1 6	Fm=6.533e16 2nd=9.061e1 5 3rd=2.918e16	Fm=3.680e16 2nd=1.229e16 3rd=1.046e16	

b.¹

a.

¹ (a) Observed SSR power, Fm=2.5 Hz. Shown are the mean SSR power values for each experimental condition, hemisphere and sensor division for the frequencies found to be significant via an *F* test (see *Methods* for details). Italicized and bolded entries indicate the comparisons found to be significant via Wilcoxon sign-rank tests between unimodal and bimodal modulation conditions. While bimodal SSR power was greater than that of the unimodal condition, not all interactions were found to be significant. The majority of significant comparisons were localized to the posterior temporal sensors for the frequency values at the modulation frequency and the theta band harmonic. (b) Identical data as in (a), Fm=3.7 Hz. Conventions are the same as described previsouly.

Fm = 3.7 Hz									
Con d	Unimodal		Φ = 0		Φ = π/2		Φ = π		
Hem	LH	RH	LH	RH	LH	RH	LH	RH	
Ant Tem	Fm=1.687e1 5 2nd=1.094e1 5 3rd=1.787e1 5	Fm=2.575e1 5 2nd=1.476e1 5 3rd=9.166e1 5	Fm=3.372e15 2nd=8.640e1 5 3rd=9.878e14	Fm=3.260e1 5 2nd=4.971e1 5 3rd=2.482e1 5	Fm=3.823e1 5 2nd=4.821e1 5 3rd=4.257e1 5	Fm=7.108e15 2nd=1.160e15 3rd=4.3526e1 5	<i>Fm=5.176e15</i> 2nd=2.792e1 5 3rd=1.620e15	<i>Fm=4.581e1</i> 5 2nd=3.791e1 5 3rd=2.560e1 5	
Pos Tem	Fm=4.701e1 5 2nd=4.938e1 5 3rd=1.170e1 6	Fm=2.575e1 5 2nd=5.001e1 5 3rd=6.959e1 5	Fm=1.384e17 2nd=2.196e1 7 3rd=5.496e15	Fm=7.631e1 6 2nd=2.451e1 6 3rd=1.156e1 6	<i>Fm</i> =7.857e1 6 2nd=3.819e1 6 3rd=1.797e1 6	Fm=1.186e17 2nd=1.134e16 3rd=9.994e15	Fm=5.400e16 2nd=4.334e1 6 3rd=1.585e16	Fm=4.149e1 6 2nd=2.452e1 6 3rd=5.488e1 5	
Occ	Fm=1.686e1 6 2nd=2.994e1 6 3rd=1.422e1 5	Fm=3.645e1 6 2nd=3.942e1 6 3rd=1.713e1 5	Fm=1.194e17 2nd=1.344e1 7 3rd=1.205e1 6	Fm=7.506e1 6 2nd=3.581e1 6 3rd=6.079e1 5	Fm=6.044e1 6 2nd=3.952e1 6 <i>3rd=1.571e1</i> 6	Fm=7.048e16 2nd=1.835e16 3rd=4.895e15	Fm=3.236e16 2nd=1.150e1 7 3rd=1.693e1 6	Fm=3.318e1 6 2nd=5.044e1 6 3rd=8.471e1 5	