

# Contributions of local speech encoding and functional connectivity to audio-visual speech integration

Bruno L. Giordano<sup>1,2\*</sup>, Robin A. A. Ince<sup>1</sup>, Joachim Gross<sup>1</sup>, Stefano Panzeri<sup>3</sup>, Philippe G. Schyns<sup>1</sup>, Christoph Kayser<sup>1\*</sup>

**\*For correspondence:**

[bruno.giordano@glasgow.ac.uk](mailto:bruno.giordano@glasgow.ac.uk) (BLG); [christoph.kayser@glasgow.ac.uk](mailto:christoph.kayser@glasgow.ac.uk) (CK)

<sup>1</sup>Institute of Neuroscience and Psychology, University of Glasgow, Glasgow, G12 8QB, UK;

<sup>2</sup>Institut de Neurosciences de la Timone UMR 7289, Aix-Marseille Université – Centre National de la Recherche Scientifique, 13005 Marseille, France; <sup>3</sup>Neural Computation Laboratory, Center for Neuroscience and Cognitive Systems, Istituto Italiano di Tecnologia, Rovereto, 38068, Italy

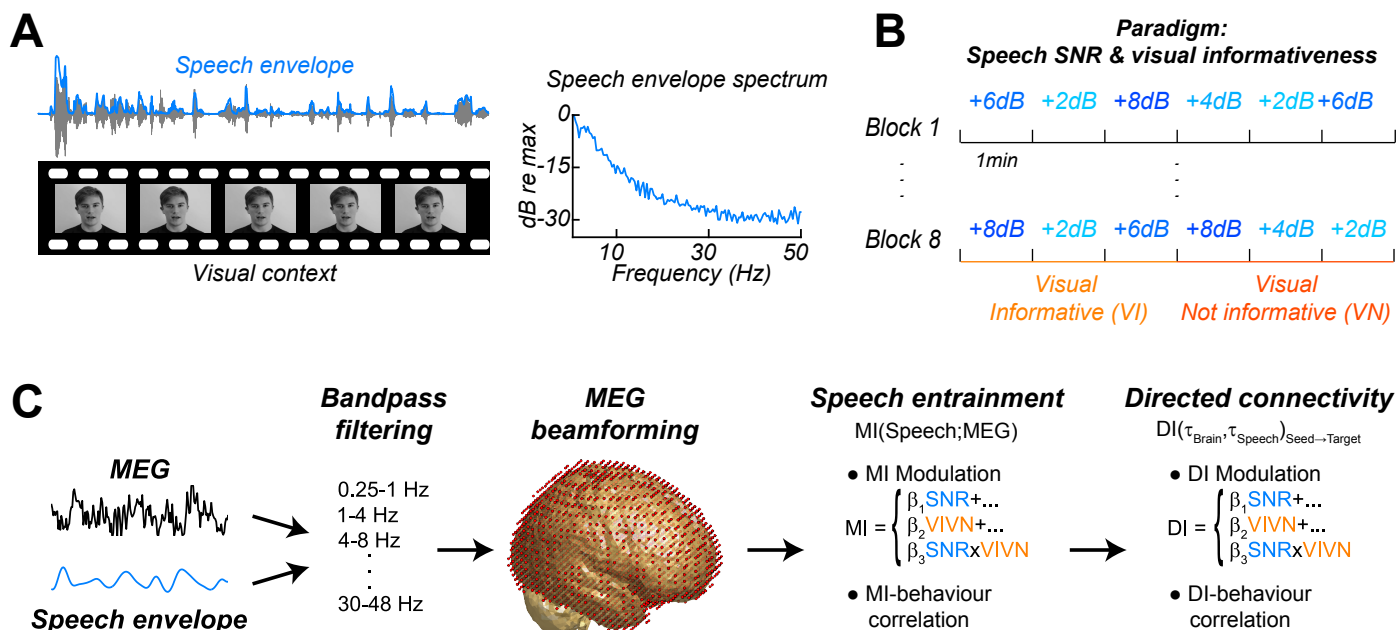
---

**Abstract** Seeing a speaker's face enhances speech intelligibility in adverse environments. We investigated the underlying network mechanisms by quantifying local speech representations and directed connectivity in MEG data obtained while human participants listened to speech of varying acoustic SNR and visual context. During high acoustic SNR speech encoding by entrained brain activity was strong in temporal and inferior frontal cortex, while during low SNR strong entrainment emerged in premotor and superior frontal cortex. These changes in local encoding were accompanied by changes in directed connectivity along the ventral stream and the auditory-premotor axis. Importantly, the behavioural benefit arising from seeing the speaker's face was not predicted by changes in local encoding but rather by enhanced functional connectivity between temporal and inferior frontal cortex. Our results demonstrate a role of auditory-motor interactions in visual speech representations and suggest that functional connectivity along the ventral pathway facilitates speech comprehension in multisensory environments.

---

## Introduction

When communicating in challenging acoustic environments we profit tremendously from visual cues arising from the speakers face. Movements of the lips, tongue or the eyes convey significant information that can boost speech intelligibility and facilitate the attentive tracking of individual speakers (*Ross et al., 2007; Sumbly and Pollack, 1954*). This multisensory benefit is strongest for continuous speech, where visual signals provide temporal markers to segment words or syllables, and provide linguistic cues (*Grant and Seitz, 1998*). Previous work has identified the synchronization of brain rhythms between interlocutors as a potential neural mechanism underlying the visual enhancement of intelligibility (*Hasson et al., 2012; Park et al., 2016; Peelle and Sommers, 2015; Pickering and Garrod, 2013; Schroeder et al., 2008*). Both acoustic and visual speech signals exhibit pseudo-rhythmic temporal structures at prosodic and syllabic rates (*Chandrasekaran et al., 2009*). These regular features can entrain rhythmic activity in the observer's brain and facilitate perception by aligning neural excitability with acoustic or visual speech features (*Giraud and Poeppel, 2012; Mesgarani and Chang, 2012; Park et al., 2016; Peelle and Davis, 2012; Schroeder et al., 2008; van Wassenhove, 2013*). While this model makes clear predictions about the visual enhancement of speech encoding in challenging environments, the network organization of multisensory speech



**Figure 1. Experimental paradigm and analysis.** (A) Stimuli consisted of 8 continuous 6 minute long audio-visual speech samples. (B) The experimental design comprised 8 conditions, defined by the factorial combination of 4 levels of speech to background signal to noise ratio (SNR = 2, 4, 6, and 8 dB) and two levels of visual informativeness (VI: Visual context Informative: video showing the narrator in synchrony with speech; VN: Visual context Not informative: video showing the narrator producing babble speech). Experimental conditions lasted 1 (SNR) or 3 (VIVN) minutes, and were presented in pseudo-randomized order. (C) Analyses were carried out on band-pass filtered speech envelope and MEG signals. The MEG data were source-projected onto a grey-matter grid (LCMV beamformer). One analysis quantified speech entrainment, i.e. the mutual information (MI) between the MEG data and the speech envelope, and the extent to which this was modulated by the experimental conditions. A second analysis quantified directed functional connectivity (DI) between seeds and the extent to which this was modulated by the experimental conditions. A final analysis assessed the correlation of either MI or DI with word-recognition performance.

41 enhancement remains unclear.

42 Previous work has implicated many brain regions in the visual enhancement of speech, including  
 43 superior temporal (Beauchamp et al., 2004; Nath and Beauchamp, 2011; Riedel et al., 2015; van At-  
 44 teveldt et al., 2004), premotor and inferior frontal cortices (Arnal et al., 2009; Evans and Davis, 2015;  
 45 Hasson et al., 2007b; Lee and Noppeney, 2011; Meister et al., 2007; Skipper et al., 2009; Wright et al.,  
 46 2003). Furthermore, some studies have shown that the visual facilitation of speech encoding may  
 47 even commence in early auditory cortices (Besle et al., 2008; Chandrasekaran et al., 2013; Ghaz-  
 48 anfar et al., 2005; Kayser et al., 2010; Lakatos et al., 2009; Zion Columbia et al., 2013). However, it  
 49 remains to be understood whether visual context shapes the encoding of speech differentially within  
 50 distinct regions of the auditory pathways, or whether the visual facilitation observed within auditory  
 51 regions is simply fed forward to upstream areas, perhaps without further modification. Hence, it  
 52 is still unclear whether the enhancement of speech-to-brain entrainment is a general mechanism  
 53 that mediates visual benefits at multiple stages along the auditory pathways.

54 Many previous studies on this question were limited by three conceptual shortcomings: first,  
 55 many have focused on generic brain activations rather than directly mapping the task-relevant sen-  
 56 sory representations (activation mapping vs. information mapping, Kriegeskorte et al., 2006), and  
 57 hence have not quantified multisensory influences on those neural representations directly shaping  
 58 behavioural performance. Second, while many studies have correlated speech-induced local brain  
 59 activity with behavioural performance, few studies have quantified directed connectivity along the  
 60 auditory pathways to ask whether perceptual benefits are better explained by changes in local en-  
 61 coding or by changes in functional connectivity. And third, most studies have neglected the contin-  
 62 uous predictive structure of speech by focusing on isolated words or syllables. However, this structure

63 may play a central role for mediating the visual benefits (*Bernstein et al., 2004; Giraud and Poeppel,*  
64 *2012; Schroeder et al., 2008; Schwartz and Savariaux, 2014*). Importantly, given that the predictive  
65 visual context interacts with acoustic signal quality to increase perceptual benefits in adverse envi-  
66 ronments (*Callan et al., 2014; Ross et al., 2007; Schwartz et al., 2004; Sumbly and Pollack, 1954*),  
67 one needs to manipulate both factors to fully address this question. Overcoming these problems,  
68 we capitalized on the statistical and conceptual power offered by naturalistic speech to study the  
69 network mechanisms that underlie the visual facilitation of speech perception.

70 Using source localized MEG activity we systematically investigated how local speech representa-  
71 tions and task-relevant directed functional connectivity along the auditory pathways change with  
72 visual context and acoustic signal quality. Specifically, we extracted neural signatures of speech  
73 representations by quantifying the mutual information between the MEG signal and the speech  
74 envelope. Furthermore, we quantified directed causal connectivity between nodes in the speech  
75 network using lagged mutual information between MEG source signals. Using linear modelling we  
76 then asked how local encoding and connectivity are affected by contextual information about the  
77 speakers face, by the acoustic signal to noise ratio, and by their interaction, and how each of these  
78 neural signatures relates to behavioural performance.

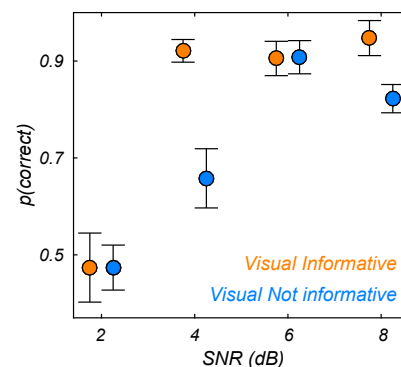
## 79 Results

80 Participants ( $n = 19$ ) were presented with continuous speech  
81 that varied in acoustic quality (signal to noise ratio, SNR) and  
82 the informativeness of the speaker's face. The visual con-  
83 text could be either informative (VI), showing the face pro-  
84 ducing the acoustic speech, or uninformative (VN), showing  
85 the same face producing nonsense babble (Fig. 1A,B). We  
86 measured brain-wide activity using MEG while participants  
87 listened to eight six-minute texts and performed a delayed  
88 word recognition task. Behavioural performance was better  
89 during high SNR and an informative visual context (Fig. 2):  
90 a repeated measures ANOVA revealed a significant effect of  
91 SNR ( $F(3,54) = 36.22$ ,  $p < 0.001$ , Huynh-Feldt corrected,  $\eta_p^2 =$   
92  $0.67$ ), and of visual context ( $F(1,18) = 18.95$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.51$ ),  
93 as well as a significant interaction ( $F(3,54) = 4.34$ ,  $p = 0.008$ ,  $\eta_p^2 =$   
94  $0.19$ ). This interaction arose from a significant visual enhance-  
95 ment for SNRs of 4 and 8 dB (paired  $T(18) \geq 3.00$ , Bonferroni  
96 corrected  $p \leq 0.032$ ;  $p > 0.95$  for other SNRs).

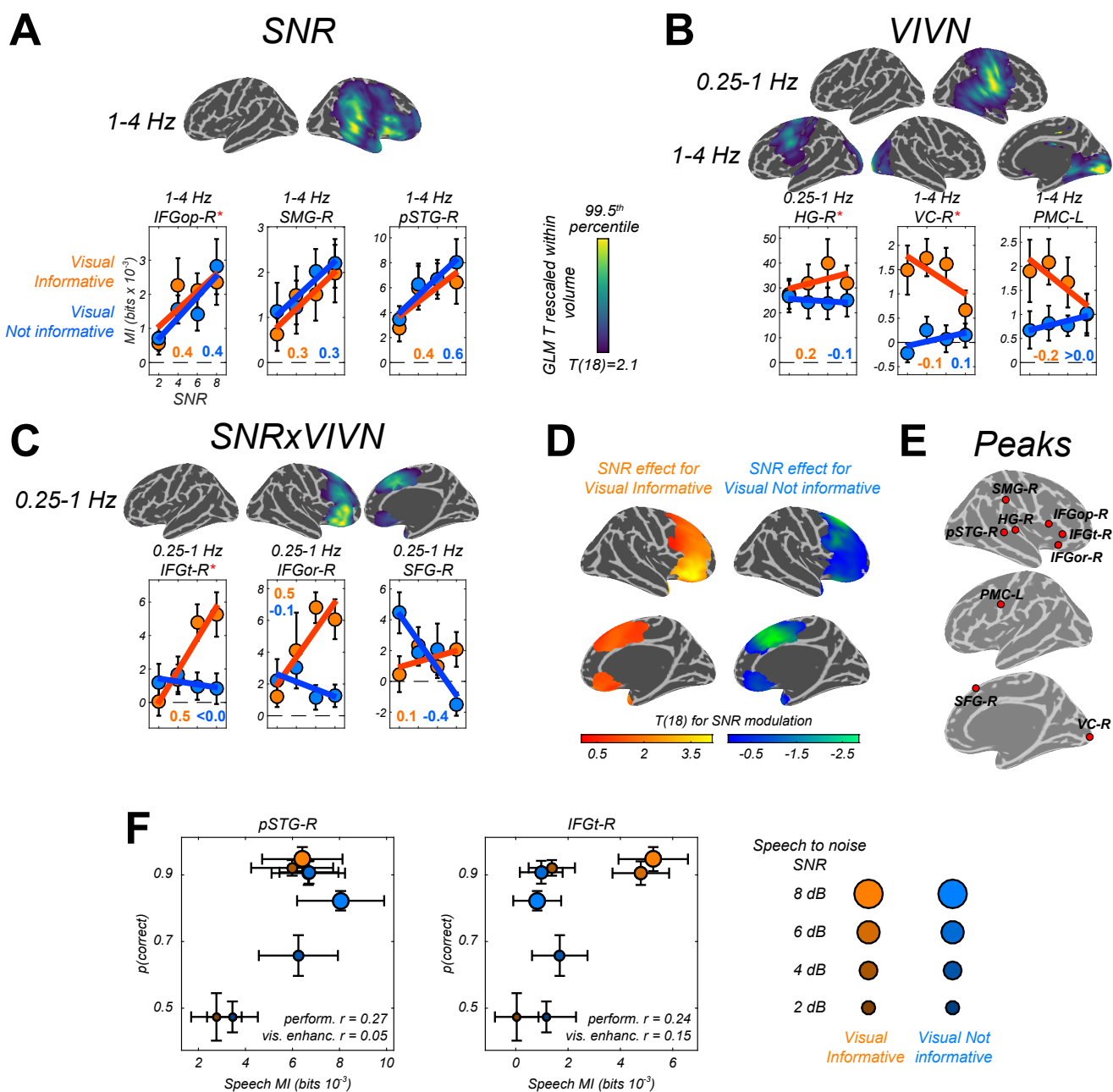
97 To study the brain activity underlying this behavioral ben-  
98 efit we analyzed source-projected MEG data using information theoretic tools to quantify the fidelity  
99 of local neural representations of the speech envelope (speech-to-brain entrainment), as well as the  
100 directed causal connectivity between relevant regions. For both, coding and connectivity, we (1)  
101 modelled the extent to which they were modulated by the experimental conditions and (2) asked  
102 whether they correlated with behavioural performance across conditions and with the visual benefit  
103 (VI-VN) across SNRs (Fig. 1C).

## 104 Widespread speech-to-brain entrainment at multiple time scales

105 Speech-to-brain entrainment was quantified by the mutual information (MI) between the MEG time  
106 course and the speech envelope (not the speech + noise mixture) in individual frequency bands  
107 (*Gross et al., 2013; Kayser et al., 2015b*, Fig. 1). At the group-level, we observed widespread signif-  
108 icant speech MI in all considered bands from 0.25 to 48 Hz (FWE = 0.05), except between 18–24  
109 Hz (Fig. S1A). Consistent with previous results (*Gross et al., 2013; Ng et al., 2013; Park et al., 2016*)  
110 speech MI was higher at low frequencies and strongest below 4 Hz (Fig. S1B). This time scale is typ-  
111 ically associated with syllabic boundaries or prosodic stress (*Giraud and Poeppel, 2012; Greenberg*



**Figure 2. Behavioural performance.** Word recognition performance for each of the experimental conditions (mean  $\pm$  SEM across participants  $n=19$ ).



**Figure 3. Modulation of speech-to-brain entrainment by acoustic SNR and visual informativeness.** Changes in speech entrainment with the experimental factors were quantified using a GLM for the condition-specific speech MI based on the effects of SNR (**A**), visual informativeness VIVN (**B**), and their interaction (SNRxVIVN) (**C**). The figures display the cortical-surface projection onto the Freesurfer template (proximity = 10 mm) of the group-level significant statistics for each GLM effect (FWE = 0.05). Graphs show the average speech MI values for each condition (mean  $\pm$  SEM), for local and global (red asterisk) of the T maps. Lines indicate the across-participant average regression model and numbers indicate the group-average standardized regression coefficient for SNR in the VI and VN conditions. (**D**) T maps illustrating the opposite SNR effects within voxels with significant SNRxVIVN effects. MI graphs for the peaks of these maps are shown in (C) (IFGor-R and SFG-R = global T peaks for SNR effects in VI and VN, respectively). (**E**) Location of global and local seeds of GLM T maps, used for the analysis of directed connectivity. (**F**) Correlation between condition-specific behavioural performance and speech MI (perform.  $r$ ) and between visual enhancement of performance and MI (vis. enhanc.  $r$ ; see inset) in pSTG-R and IFGt-R. error-bars =  $\pm$  SEM. See also Tables 1 and 3.

112 **et al., 2003**). Indeed, the average syllabic rate was 212 syllables per minute in the present material,  
 113 corresponding to about 3.5 Hz. Across frequencies, MI was strongest in bilateral auditory cortex and  
 114 more extended within the right hemisphere (Fig. S1B). Peak MI values were significantly higher in the

**Table 1. Condition effects on speech MI.** The table lists global and local peaks in the GLM T-maps. Anatomical labels and Brodmann areas are based on the AAL and Talairach atlases.  $\beta$  = standardized regression coefficient; SEM = standard error of the participant average.

Anatomical label	Brodmann area	MNI coordinates			GLM effect	Frequency Band	T (18)	$\beta$ (SEM)
HG-R	42	64	-20	12	VIVN	0.25-1 Hz	4.62	0.10(0.15)
pSTG-R	22	48	-30	8	SNR	1-4 Hz	4.46	0.48(0.08)
SMG-R	40	58	-30	38	SNR	1-4 Hz	3.9	0.29(0.09)
PMC-L	6	-54	0	32	VIVN	1-4 Hz	3.81	0.62(0.20)
IFGt-R	46	42	34	2	SNRxVIVN	0.25-1 Hz	3.62	0.66(0.15)
IFGop-R	47	50	18	2	SNR	1-4 Hz	4.94	0.36(0.08)
IFGor-R	47	30	26	-16	SNR in VI	0.25-1 Hz	5.04	0.48(0.09)
SFG-R	6	12	30	58	SNR in VN	0.25-1 Hz	-3.54	-0.44(0.10)
VC-R	17/18	18	-102	-4	VIVN	1-4 Hz	6.01	0.72(0.15)

115 right compared to the left hemisphere at frequencies below 12 Hz (paired t-tests;  $T(18) \geq 3.1$ ,  $p \leq 0.043$   
 116 Bonferroni corrected), and did not differ at higher frequencies ( $T(18) \leq 2.78$ ,  $p \geq 0.08$ ). Importantly,  
 117 we observed significant speech-to-brain entrainment not only within temporal cortices but across  
 118 multiple regions in the occipital, frontal and parietal lobes, consistent with the notion that speech  
 119 information is represented also within motor and frontal regions (*Bornkessel-Schlesewsky et al.,*  
 120 *2015; Du et al., 2014; Skipper et al., 2009*).

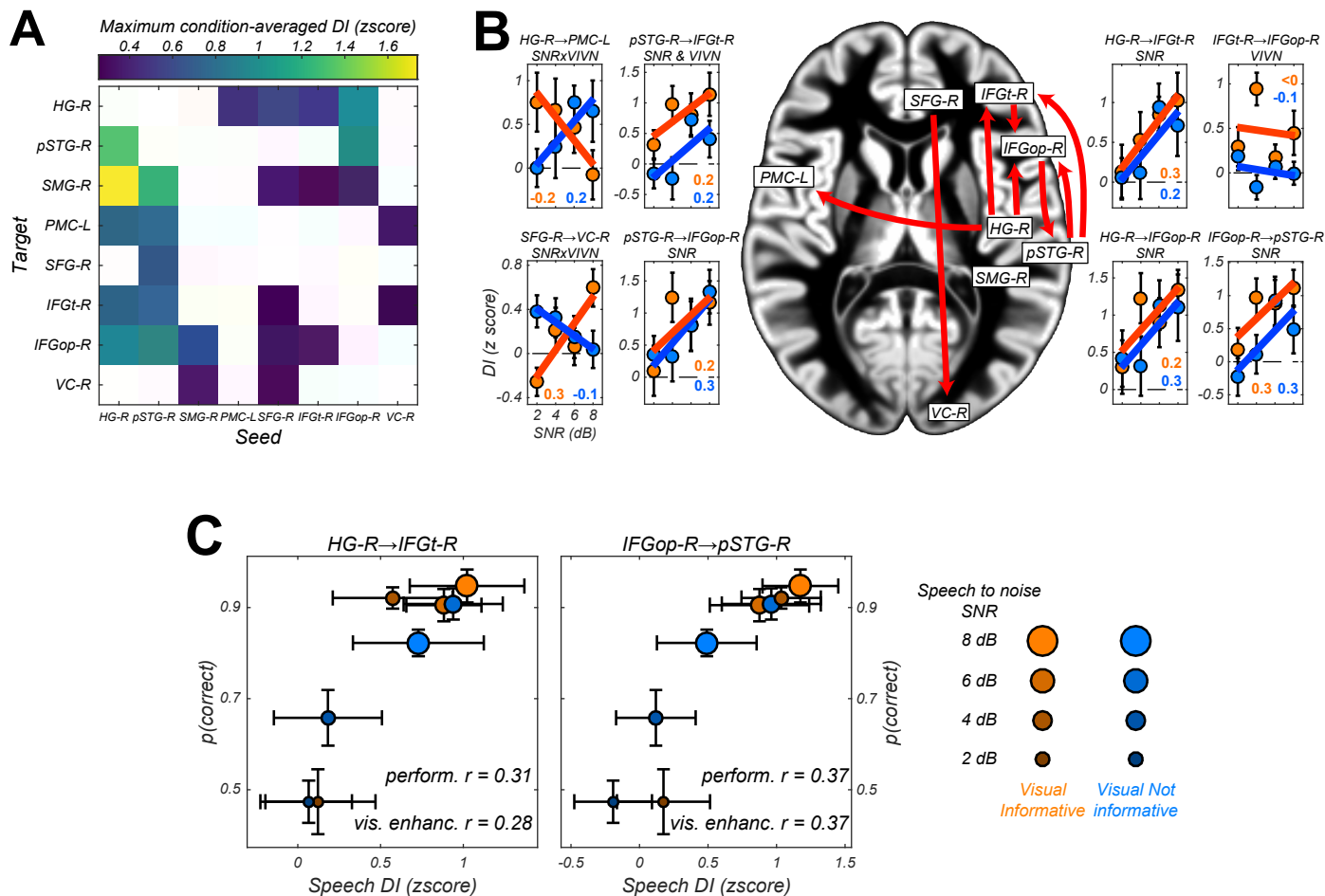
### 121 **Speech entrainment is modulated by SNR within and beyond auditory cortex**

122 To determine the regions where acoustic signal quality and visual context affect the encoding of  
 123 speech we modelled the condition-specific speech MI values based on effects of acoustic signal  
 124 quality (SNR), visual informativeness (VIVN), and their interaction (SNRxVIVN). Random-effects signif-  
 125 icance was tested using a permutation procedure and cluster enhancement, correcting for multiple  
 126 comparisons along all relevant dimensions. Effects of experimental factors emerged in multiple re-  
 127 gions at frequencies below 4 Hz (Fig. 3). Increasing the acoustic signal quality (SNR; Fig. 3A) resulted  
 128 in stronger speech MI in the right auditory cortex (1-4 Hz; local peak T statistic = 4.46 in posterior supe-  
 129 rior temporal gyrus; pSTG-R; Table 1), right parietal cortex (local peak T = 3.90 in supramarginal gyrus;  
 130 SMG-R), and right dorso-ventral frontal cortex (IFGop-R; global peak T = 4.94). We also observed sig-  
 131 nificant positive SNR effects within the right temporo-parietal and occipital cortex at 12-18 Hz (local  
 132 peak right lingual gyrus, T = 5.12). However, inspection of the participant-specific data suggested that  
 133 this effect was not reliable (for only 58% of participants showed an speech MI increase with SNR, as  
 134 opposed to a minimum of 84% for the other SNR effects), possibly because the comparatively lower  
 135 power of speech envelope fluctuations at higher frequencies (c.f. Fig. 1A), and hence this effect is  
 136 not discussed further.

### 137 **Visual context reveals distinct strategies for handling speech in noise in premotor, su- 138 perior and inferior frontal cortex**

139 Contrasting informative and un-informative visual contexts revealed stronger speech MI when see-  
 140 ing the speakers face (VI) at frequencies below 4 Hz in both hemispheres (Fig. 3B): the right temporo-  
 141 parietal cortex (0.25-1 Hz; HG; T = 4.62), bilateral occipital cortex (1-4 Hz; right visual cortex; VC-R; global  
 142 T peak = 6.01) and left premotor cortex (1-4 Hz; PMC-L; local T peak = 3.81). Interestingly, the condition-  
 143 specific pattern of MI for both VC-L and PMC-L were characterized by an increase in speech MI with  
 144 decreasing SNR during the VI condition, pointing to a stronger visual enhancement during more  
 145 adverse listening conditions.

146 Since visual benefits for perception emerge mostly when acoustic signals are degraded (Fig. 2,  
 147 *Ross et al., 2007; Sumbly and Pollack, 1954*), the interaction of acoustic and visual factors provides



**Figure 4. Directed causal connectivity within the speech-entrained network.** Directed connectivity between seeds of interest (c.f. Fig. 3E) was quantified using Directed Information (DI). **(A)** Maximum significant condition-average DI across lags (FWE = 0.05 across lags; white = no significant DI). **(B)** Significant condition effects (GLM for SNR, VIVN or their interaction) on DI (FWE = 0.05 across speech/brain lags and seed/target pairs). Bar graphs display condition-specific DI values for each significant GLM effect along with the across-participants average regression model (lines). Numbers indicate the group-average standardized betas for SNR in the VI and VN conditions, averaged across lags associated with a significant GLM effect. **(C)** Correlation between behavioural performance and condition-specific DI (*perform. r*) and between visual enhancement of performance and DI (*vis. enhanc. r*) from HG-R to IFGt-R and from IFGop-R to pSTG-R. error-bars =  $\pm$  SEM. See also Tables 2-3 and Fig. S2.

148 a crucial test for audio-visual integration. We found significant interactions in the 0.25-1 Hz band in  
 149 the right dorso-ventral frontal lobe, which peaked in the pars triangularis (IFGt-R;  $T = 3.62$ ; Fig. 3C).  
 150 Importantly, investigating the SNR effect at these voxels revealed two distinct strategies for handling  
 151 speech in noise dependent on visual context (Fig. 3D): During VI speech MI increased with SNR in  
 152 ventral frontal cortex (peak  $T$  for SNR in pars orbitalis; IFGop-R;  $T = 5.04$ ), while in dorsal frontal cortex  
 153 speech MI was strongest at low SNRs during VN (peak  $T$  in superior frontal gyrus; SFG-R;  $T = -3.54$ ). This  
 154 demonstrates distinct functional roles of ventral and dorsal prefrontal regions in speech encoding  
 155 and reveals a unique role of superior frontal cortex for enhancing speech representations in a poorly  
 156 informative context, such as the absence of visual information in conjunction with poor acoustic  
 157 signals.

### 158 Directed causal connectivity within the speech network

159 The diversity of the patterns of speech entrainment in temporal, premotor and inferior frontal re-  
 160 gions across conditions could arise from the individual encoding properties of each region, or from  
 161 changes in functional connectivity between regions with conditions. To directly test this, we quanti-

**Table 2. Analysis of directed connectivity (DI).** The table lists connections with significant condition-averaged DI, and condition effects on DI. SEM = standard error of participant average;  $\beta$  = standardized regression coefficients. T(18) = maximum T statistic within significance mask.

Seed	Target	DI T(18)	Condition effects (GLM)		
			Effect	T(18)	$\beta$ (SEM)
HG-R	PMC-L	3.38	SNRxVIVN	-3.01	-0.14(0.05)
HG-R	IFGt-R	3.03	SNR	3.32	0.18(0.05)
HG-R	IFGopR	4.54	SNR	3.19	0.18(0.05)
pSTG-R	IFGt-R	3.39	SNR	3.91	0.22(0.06)
			VIVN	4.57	0.59(0.22)
pSTG-R	IFGopR	4.12	SNR	3.31	0.20(0.06)
SFG-R	VC-R	4.4	SNRxVIVN	3.69	0.12(0.03)
IFGt-R	IFGopR	3.76	VIVN	3.56	0.31(0.17)
IFGopR	pSTG-R	4.16	SNR	4.65	0.17(0.04)

162 fied the directed causal connectivity between regions of interest extracted from the above analysis  
 163 (Fig. 3E). To this end we used Directed Information (DI), also known as Transfer Entropy, an infor-  
 164 mation theoretic measure of Wiener-Granger causality (Massey, 1990; Schreiber, 2000). We took  
 165 advantage of previous work that made this measure statistically robust when applied to neural data  
 166 (Besserve et al., 2015; Ince et al., 2016a).

167 We observed significant condition-averaged DI between multiple nodes of the speech network  
 168 (FWE = 0.05; Fig. 4A and Fig. S2A). This included among others the feed-forward pathways of the ven-  
 169 tral and dorsal auditory streams, such as from auditory cortex (HG-R) and superior temporal regions  
 170 (pSTG-R) to premotor (PMC-L) and to inferior frontal regions (IFGt-R, IFGop-R), from right parietal  
 171 cortex (SMG-R) to premotor cortex (PMC-L), as well as feed-back connections from premotor and  
 172 inferior frontal regions to temporal regions. In addition, we also observed significant connectivity  
 173 between frontal (SFG-R) and visual cortex (VC).

174 We then asked whether and where connectivity changed with experimental conditions (Fig. 4B,  
 175 Table 2 and Fig. S2B). Within the right ventral stream feed-forward connectivity from the tempo-  
 176 ral lobe (HG-R, pSTG-R) to frontal cortex (IFGt-R, IFGop-R) was enhanced during high acoustic SNR  
 177 (FWE = 0.05;  $T(18) \geq 3.1$ ). More interestingly, this connectivity was further enhanced in the presence of  
 178 an informative visual context (pSTG-R  $\rightarrow$  IFGt-R, positive SNRxVIVN interaction,  $T = 4.57$ ), demonstrat-  
 179 ing a direct influence of visual context on the propagation of speech information along the ventral  
 180 stream. Interactions of acoustic and visual context on connectivity were also found from auditory  
 181 (HG-R) to premotor cortex (PMC-L, negative interaction;  $T = -3.01$ ). Here connectivity increased with  
 182 increasing SNR in the absence of visual information and increased with decreasing SNR during an  
 183 informative context, suggesting that visual information changes the qualitative nature of auditory-  
 184 motor interactions. An opposite interaction was observed between the frontal lobe and visual cortex  
 185 (SFG-R  $\rightarrow$  VC-R,  $T = 4.40$ ). Finally, we found that feed-back connectivity along the ventral pathway  
 186 was significantly stronger during high SNRs (IFGt-R  $\rightarrow$  pSTG-R;  $T = 4.16$ ).

### 187 Do Speech entrainment or connectivity shape behavioural performance?

188 We performed two additional analyses to test whether and where changes in the local represen-  
 189 tation of speech information (speech-MI) or directed connectivity (DI) contribute to explaining the  
 190 behavioural benefits (Fig. 2). First, we asked where speech-MI/DI relates to performance changes  
 191 across all experimental conditions (incl. changes in SNR). This revealed a significant correlation be-  
 192 tween condition-specific word-recognition performance and the strength of speech MI in pSTG-R  
 193 and IFGt - R ( $r \geq 0.28$ ; FWE = 0.05; Table 3 and Fig. 3F), suggesting that stronger entrainment in

**Table 3. Association of behavioural performance with speech entrainment and connectivity.** Performance: T statistic and average of participant-specific correlation (SEM) between behavioural performance and speech MI / DI. Visual enhancement: correlation between SNR-specific behavioural benefit (VI-VN) and the respective difference in speech-MI or DI. \* FWE = 0.05 corrected for multiple comparisons.

<b>Speech MI</b>					
		<b>Performance</b>		<b>Visual enhancement</b>	
		<b>T(18)</b>	<b>r(SEM)</b>	<b>T(18)</b>	<b>r(SEM)</b>
	HG-R	1.27	0.12(0.08)	0.21	0.04(0.12)
	pSTG-R	3.43 *	0.27(0.07)	0.53	0.05(0.10)
	SMG-R	2.35	0.19(0.08)	-0.39	-0.02(0.10)
	PMC-L	0.47	0.04(0.07)	0.13	0.03(0.12)
	SFG-R	-0.47	-0.03(0.07)	1.61	0.17(0.11)
	IFGt-R	3.09 *	0.24(0.08)	1.25	0.15(0.12)
	IFGopR	2.38	0.20(0.08)	-0.25	-0.01(0.12)
	VC-R	1.55	0.14(0.09)	-0.82	-0.16(0.11)

<b>Directed connectivity</b>					
		<b>Performance</b>		<b>Visual enhancement</b>	
<b>Seed</b>	<b>Target</b>	<b>T(18)</b>	<b>r(SEM)</b>	<b>T(18)</b>	<b>r(SEM)</b>
HG-R	IFGt-R	4.83 *	0.31(0.07)	2.55 *	0.28(0.11)
HG-R	IFGopR	3.19 *	0.24(0.07)	1.86	0.31(0.17)
HG-R	PMC-L	0.90	0.06(0.06)	-0.07	-0.01(0.14)
pSTG-R	IFGt-R	4.28 *	0.27(0.06)	1.28	0.16(0.12)
pSTG-R	IFGopR	3.59 *	0.29(0.08)	1.82	0.32(0.17)
IFGt-R	IFGopR	1.11	0.08(0.07)	2.27	0.33(0.14)
IFGopR	pSTG-R	4.51 *	0.37(0.08)	2.55 *	0.37(0.15)
SFG-R	VC-R	-0.04	0.00(0.08)	0.90	0.17(0.18)

194 the ventral stream facilitates comprehension. This hypothesis was further corroborated by a signif-  
 195 icant correlation of connectivity along the ventral stream with behavioural performance, both in  
 196 feed-forward (HG-R → IFGt-R; pSTG-R → IFGt-R/IFGop-R;  $r \geq 0.27$ , Table 3) and feed-back directions  
 197 (IFGop-R → pSTG-R;  $r=0.37$ ). The enhanced quality of speech perception during favourable listening  
 198 conditions hence results from enhanced speech encoding and the supporting network connections  
 199 along the temporal-frontal axis.

200 Second, we asked whether and where the improvement in behavioural performance with an in-  
 201 formative visual context (VI-VN) correlates with an enhancement in speech encoding or connectivity.  
 202 This revealed no significant correlation between the visual enhancement of local speech representa-  
 203 tions and perceptual benefits (all  $p > 0.05$ ). However, both feed-forward (HG-R → IFGt-R;  $r = 0.28$ ,  $p <$   
 204  $0.05$ ; Fig. 4C) and feed-back connections (IFGop-R → pSTG-R;  $r = 0.37$ ) along the ventral stream were  
 205 significantly enhanced during an informative visual context, suggesting that changes in functional  
 206 connectivity contribute significantly to shaping speech intelligibility.

## 207 Discussion

208 The present study provides a comprehensive picture of how acoustic signal quality and visual con-  
 209 text interact to shape the encoding of speech information and the directed functional connectivity  
 210 along speech-sensitive cortex. Our results reveal a dominance of feed-forward pathways from au-  
 211 ditory regions to inferior frontal cortex under favourable conditions, such as during high SNR. We



212 also demonstrate the visual enhancement of speech encoding in auditory and premotor cortex, as  
213 well as non-trivial interactions of acoustic quality and visual context in superior and inferior frontal  
214 regions. These patterns of local encoding were accompanied by changes in directed connectivity  
215 along the ventral pathway and from auditory to premotor cortex. Yet, the behavioural benefit arising  
216 from seeing the speaker's face was not related to any site-specific visual enhancement of local  
217 speech encoding. Rather, changes in directed functional connectivity along the ventral stream were  
218 predictive of the behavioural benefit.

### 219 **Entrained speech representations in temporal, parietal and frontal lobes**

220 We observed functionally distinct patterns of speech-to-brain entrainment along the auditory path-  
221 ways. Previous studies on speech entrainment largely focused on the auditory cortex, where en-  
222 trainment is strongest (*Ding and Simon, 2013; Gross et al., 2013; Keitel et al., 2017; Mesgarani and*  
223 *Chang, 2012*). This was in part due to the difficulty to separate distinct processes reflecting entrain-  
224 ment when contrasting only few experimental conditions (e.g. forward and reversed speech, *Ding*  
225 *and Simon, 2012; Gross et al., 2013*). Based on the susceptibility to changes in acoustic signal qual-  
226 ity and visual context we here establish entrainment as a ubiquitous mechanism reflecting distinct  
227 speech representations along auditory pathways.

228 Speech entrainment was reduced with decreasing acoustic SNR in temporal, parietal and ven-  
229 tral prefrontal cortex, directly reflecting the reduction in behavioural performance in challenging  
230 environments. In contrast, entrainment was enhanced during low SNR in superior frontal and pre-  
231 motor cortex. While there is strong support for a role of frontal and premotor regions in speech  
232 analysis (*Du et al., 2014; Evans and Davis, 2015; Heim et al., 2008; Meister et al., 2007; Morillon*  
233 *et al., 2015; Rauschecker and Scott, 2009; Skipper et al., 2009; Wild et al., 2012*), most evidence  
234 comes from stimulus-evoked activity rather than signatures of neural speech encoding. We directly  
235 demonstrate the specific enhancement of frontal (PMC, SFG) speech representations during chal-  
236 lenging conditions. This enhancement is not directly inherited from the temporal lobe, as temporal  
237 regions exhibited either no visual facilitation (pSTG) or visual facilitation without an interaction with  
238 SNR (HG).

239 The effects of experimental conditions dominated on the right hemisphere. Such a right domi-  
240 nance of speech entrainment is in agreement with previous studies (*Bourguignon et al., 2013; Fonte-*  
241 *neau et al., 2015; Gross et al., 2013; Vander Ghinst et al., 2016*) and with the hypothesis that right tem-  
242 poral regions extract acoustic information predominantly on the syllabic and prosodic time scales  
243 (*Giraud and Poeppel, 2012; Poeppel, 2003*), exactly those time scales where speech-to-brain entrain-  
244 ment is strongest in the present and previous data (*Gross et al., 2013; Keitel et al., 2017*).

### 245 **Multisensory enhancement of speech encoding in the frontal lobe**

246 Visual information from the speakers face provides multiple cues that enhance intelligibility. In sup-  
247 port of a behavioural multisensory benefit we found stronger entrainment during an informative  
248 visual context in multiple bilateral regions. First, we replicated the visual enhancement of audi-  
249 tory cortical representations (HG, *Besle et al., 2008; Kayser et al., 2010; Zion Golumbic et al., 2013*).  
250 Second, visual enhancement of an acoustic speech representation was also visible in early visual  
251 areas, as suggested by prior studies (*Nath and Beauchamp, 2011; Schepers et al., 2015*). While we  
252 can't rule out that this effect is in part mediated by the correlations between acoustic and visual  
253 speech cues, we found that the visual enhancement was strongest when SNR was low and hence is  
254 better explained by top-down influences (*Vetter et al., 2014*). Third, speech representations in ven-  
255 tral prefrontal cortex were selectively involved during highly reliable multisensory conditions and  
256 were reduced in the absence of the speakers face. These findings are in line with suggestions that  
257 the IFG facilitates comprehension (*Alho et al., 2014; Evans and Davis, 2015; Hasson et al., 2007b;*  
258 *Hickok and Poeppel, 2007*) and implements multisensory processes (*Callan et al., 2014, 2003; Lee*  
259 *and Noppeney, 2011*), possibly by providing amodal phonological, syntactic and semantic processes  
260 (*Clos et al., 2014; Ferstl et al., 2008; McGettigan et al., 2012*). Previous studies often reported en-

261 hanced IFG response amplitudes under challenging conditions (*Guediche et al., 2013*). In contrast,  
262 by quantifying the fidelity of speech representations we here show that these are generally stronger  
263 during favourable SNRs. This discrepancy is not necessarily surprising, if one assumes that IFG rep-  
264 resentations are derived from those in the temporal lobe, which are also more reliable during high  
265 SNRs. Noteworthy, however, is the finding that representations within ventral IFG are selectively  
266 stronger during an informative visual context. We thereby directly confirm the hypothesis that IFG  
267 speech encoding is enhanced by visual context.

268 Furthermore, we demonstrate the visual enhancement of speech representations in premotor  
269 regions, which could implement the mapping of audio-visual speech features onto articulatory rep-  
270 resentations (*Meister et al., 2007; Morillon et al., 2015; Fernández et al., 2015; Skipper et al., 2009;*  
271 *Wilson et al., 2004*). We show that that this enhancement is inversely related to acoustic signal qual-  
272 ity. While this observation is in agreement with the notion that perceptual benefits are strongest  
273 under adverse conditions (*Ross et al., 2007; Sumbly and Pollack, 1954*), there was no significant cor-  
274 relation between the visual enhancement of premotor encoding and behavioural performance. Our  
275 results thereby deviate from previous work that has suggested a driving role of premotor regions in  
276 shaping intelligibility (*Alho et al., 2014; Osnes et al., 2011*), and we rather support a modulatory influ-  
277 ence of auditory-motor interactions (*Alho et al., 2014; Callan et al., 2004; Hickok and Poeppel, 2007;*  
278 *Krieger-Redwood et al., 2013; Morillon et al., 2015*). For example, in a study quantifying dynamic rep-  
279 resentations of visual speech signals (lip movements) we recently found that left premotor activity  
280 was significantly predictive of behavioural performance (*Park et al., 2016*). One explanation for this  
281 discrepancy may be presence of a memory component in our behavioural task, which may engage  
282 other brain regions (e.g. IFG) more than other tasks. Given that the premotor effects were restricted  
283 to the theta band, which is associated with syllabic (> 3 Hz) rather than intonational (< 1 Hz) struc-  
284 ture (*Giraud and Poeppel, 2012; Greenberg et al., 2003*), our results also suggest this region carries  
285 syllabic rather than prosodic representations (*Du et al., 2014; Heim et al., 2008; Krieger-Redwood*  
286 *et al., 2013; Osnes et al., 2011*).

287 Finally, our results highlight an interesting role of the superior frontal gyrus, where entrainment  
288 was strongest when sensory information was most impoverished (low SNR, visual not informative)  
289 or when the speakers face was combined with clear speech (high SNR, visual informative). Super-  
290 ior frontal cortex has been implied in high level inference processes underlying comprehension,  
291 sentence level integration or the exchange with memory (*Ferstl et al., 2008; Hasson et al., 2007a;*  
292 *Yarkoni et al., 2008*) and is sometimes considered part of the broader semantic network (*Binder*  
293 *et al., 2009; Gow and Olson, 2016; Price, 2012*). Our data show that the SFG plays a critical role for  
294 speech encoding under challenging conditions at the supra-syllabic time scale, possibly by medi-  
295 ating sentence-level integration during low SNRs or the comparison of visual prosody with acoustic  
296 inputs in multisensory contexts.

### 297 **Multisensory behavioural benefits arise from distributed network mechanisms**

298 To understand whether the condition-specific patterns of local speech representations emerge within  
299 each region, or whether they are possibly established by network interactions we investigated the  
300 directed functional connectivity between regions of interest. While many studies have assessed the  
301 connectivity between auditory regions (e.g. *Abrams et al., 2013; Chu et al., 2013; Fonteneau et al.,*  
302 *2015; Park et al., 2015*), few have quantified the behavioural relevance of these connections (*Alho*  
303 *et al., 2014*).

304 We observed significant intra-hemispheric connectivity between right temporal, parietal and  
305 frontal regions, in line with the transmission of speech information from auditory cortices along the  
306 auditory pathways (*Bornkessel-Schlesewsky et al., 2015; Hickok, 2012; Poeppel, 2014*). Support-  
307 ing the idea that acoustic representations are progressively transformed along these pathways we  
308 found that the condition-specific patterns of functional connectivity differed systematically along  
309 the ventral and dorsal streams. While connectivity along the ventral stream was predictive of be-  
310 havioural performance and strongest during favourable listening conditions, the inter-hemispheric

311 connectivity to left premotor cortex was strongest during adverse multisensory conditions. Our re-  
312 sults therefore suggest that premotor representations are informed by auditory regions (HG, pSTG)  
313 rather than being driven by the frontal lobe, an interpretation that is supported by previous work  
314 (*Alho et al., 2014; Gow and Olson, 2016; Osnes et al., 2011*).

315 Across changes in visual context and acoustic SNR behavioural performance was supported  
316 both by an enhancement of speech representations along multiple regions of the ventral pathway  
317 and increases in their functional connectivity. These increases in functional connectivity emerged  
318 both along feed-forward and feed-back directions between temporal and inferior frontal regions,  
319 and were strongest (in effect size) along the feed-back route. This underlines the hypothesis that  
320 recurrent processing, rather than a simple feed-forward sweep, is central to speech intelligibility  
321 (*Bornkessel-Schlesewsky et al., 2015; Hickok, 2012; Poeppel, 2014*). Central to the scope of the  
322 present study, however, we found that no single region-specific effect could explain the visual be-  
323 havioural benefit. Rather, the benefit arising from seeing the speakers face was significantly corre-  
324 lated with the enhancement of functional connectivity along the ventral stream (HG → IFG → pSTG).  
325 Our results hence point to a distributed origin of the visual enhancement of speech intelligibility.  
326 As proposed a decade ago (*Besle et al., 2008; Ghazanfar et al., 2005; Ghazanfar and Schroeder,*  
327 *2006; Kayser et al., 2010; Zion Columbic et al., 2013*) this visual enhancement involves early audi-  
328 tory cortices, but as we show here, the behavioural benefit also relies on the recurrent transformation  
329 of speech representations between temporal and frontal regions. One interpretation of this in the  
330 context of predictive coding models is that an informative visual context facilitates the correction  
331 of prior predictions about the expected stimulus by incoming sensory evidence, which would be  
332 visible both in feed-forward and feed-back connectivity (*Arnal and Giraud, 2012; Bastos et al., 2012*).

333 Our results provide a network view on the dynamic speech representations in multisensory en-  
334 vironments. While premotor and superior frontal regions are specifically engaged in the most chal-  
335 lenging environments the visual enhancement of comprehension at intermediate SNRs is mediated  
336 by interactions of the core speech regions along the ventral pathway. Such a distributed neural ori-  
337 gin of multisensory benefits is in line with the notion of a hierarchical organization of multisensory  
338 processing in the brain (*Lee and Noppeney, 2011; Rohe and Noppeney, 2015*), and the idea that  
339 comprehension is shaped by network connectivity more than the engagement of particular brain  
340 regions (*Abrams et al., 2013*).

## 341 **Materials and methods**

342 Nineteen right handed healthy adults (10 females; age from 18 to 37) participated in this study. All  
343 participants were tested for normal hearing, were briefed about the nature and goal of this study,  
344 and received financial compensation for their participation. The study was conducted in accordance  
345 with the Declaration of Helsinki and was approved by the local ethics committee (College of Science  
346 and Engineering, University of Glasgow). Written informed consent was obtained from all partici-  
347 pants.

## 348 **Stimulus material**

349 The stimulus material consisted of audio-visual recordings based on text transcripts taken from pub-  
350 licly available TED talks also used in a previous study (*Kayser et al., 2015b*, Fig. 1A). Acoustic (44.1 kHz  
351 sampling rate) and video recordings (25 Hz frame rate, 1920 by 1080 pixels) were obtained while  
352 a trained male native English speaker narrated these texts (*Kayser et al., 2015a*). The root mean  
353 square (RMS) intensity of each audio recording was normalized using 6 s sliding windows to ensure  
354 a constant average intensity. Across the eight texts the average speech rate was 160 words (range  
355 138–177) per minute, and the syllabic rate was 212 syllables (range 192–226) per minute.

## 356 **Experimental design and stimulus presentation**

357 We presented each of the eight texts as continuous 6 minute sample, while manipulating the acous-  
358 tic quality and the visual relevance in a block design within each text (Fig. 1B). The visual relevance

359 was manipulated by either presenting the video matching the respective speech (visual informative,  
360 VI) or presenting a 3 s babble sequence that was repeated continuously (visual not informative, VN),  
361 and which started and ended with the mouth closed to avoid transients. The signal to noise ratio  
362 (SNR) of the acoustic speech was manipulated by presenting the speech on background cacophony  
363 of natural sounds and scaling the relative intensity of the speech while keeping the intensity of the  
364 background fixed. We used relative SNR values of +8, +6, +4 and +2 dB RMS intensity levels. The acous-  
365 tic background consisted of a cacophony of naturalistic sounds, created by randomly superimposing  
366 various naturalistic sounds from a larger database (using about 40 sounds at each moment in time,  
367 **Kayser et al., 2016**). This resulted in a total of 8 conditions (four SNR levels; visual informative or irrel-  
368 evant) that were introduced in a block design (Fig. 1B). The SNR changed from minute to minute in  
369 a pseudo-random manner (12 one minute blocks per SNR level). Visual relevance was manipulated  
370 within 3 minute sub-blocks. Texts were presented with self-paced pauses. Subjects performed a de-  
371 layed comprehension tasks after each block, whereby they had to indicate whether a specific word  
372 (noun) was mentioned in the previous text (6 words per text) or not (6 words per text) in a two alter-  
373 native forced choice task. The words chosen from the presented text were randomly selected and  
374 covered all eight conditions. The average performance was  $73 \pm 2\%$  correct (mean and SEM across  
375 subjects), showing that they indeed paid attention to the stimulus. Behavioural performance was  
376 averaged within each condition, and analysed using a repeated measures ANOVA, with SNR and  
377 VIVN as within-subject factors. The stimulus presentation was controlled using the Psychophysics  
378 toolbox in Matlab (**Brainard, 1997**). Acoustic stimuli were presented using an Etymotic ER-30 tube-  
379 phone (tube length = 4 m) at 44.1 kHz sampling rate and an average intensity of 65 dB RMS level,  
380 calibrated separately for each ear. Visual stimuli were presented in grey-scale and projected onto a  
381 translucent screen at  $1280 \times 720$  pixels at 25 fps covering a field of view of  $41 \times 33$  degrees.

### 382 **Pre-processing of the speech envelope**

383 We extracted the envelope of the speech signal (not the speech plus background mixture) by com-  
384 puting the wide-band envelope at 150 Hz temporal resolution as in previous work (**Chandrasekaran**  
385 **et al., 2009; Kayser et al., 2015b**). The speech signal was filtered (4<sup>th</sup> order Butterworth filter; forward  
386 and reverse) into six frequency bands (100 Hz–4 kHz) spaced to cover equal widths on the cochlear  
387 map. The wide-band envelope was defined as the average of the Hilbert envelopes of these band-  
388 limited signals (c.f. Fig. 1A).

### 389 **MEG data collection**

390 MEG recordings were acquired with a 248-magnetometers whole-head MEG system (MAGNES 3600  
391 WH, 4-D Neuroimaging) at a sampling rate of 1017.25 Hz. Participants were seated upright. The  
392 position of five coils, marking fiducial landmarks on the head of the participants, was acquired at  
393 the beginning and at the end of each block. Across blocks, and participants, the maximum change  
394 in their position was 3.6 mm, on average (STD = 1.2 mm).

### 395 **MEG pre-processing**

396 Analyses were carried out in Matlab using the Fieldtrip toolbox (**Oostenveld et al., 2010**), SPM12,  
397 and code for the computation of information-theoretic measures (**Ince et al., 2016a**). Block-specific  
398 data were pre-processed separately. Infrequent SQUID jumps (observed in 1.5% of the channels, on  
399 average) were repaired using piecewise cubic polynomial interpolation. Environmental magnetic  
400 noise was removed using regression based on principal components of reference channels. Both  
401 the MEG and reference data were filtered using a forward-reverse 70 Hz FIR low-pass (-40 dB at  
402 72.5 Hz); a 0.2 Hz elliptic high-pass (-40 dB at 0.1 Hz); and a 50 Hz FIR notch filter (-40 dB at  $50 \pm$   
403 1Hz). Across participants and blocks, 7 MEG channels were discarded as they exhibited a frequency  
404 spectrum deviating consistently from the median spectrum (shared variance < 25%). For analysis  
405 signals were resampled to 150 Hz, high-pass filtered at 0.2 Hz (forward-reverse elliptic filter). ECG and

406 EOG artefacts were removed using ICA in fieldtrip (runica on 40 principal components), and were  
407 identified based on the time course and topography of IC components (**Hipp and Siegel, 2013**).

### 408 **Structural data and source localization**

409 High resolution anatomical MRI scans were acquired for each participant (voxel size = 1 mm<sup>3</sup>) and co-  
410 registered to the MEG data using a semi-automated procedure. Anatomicals were segmented into  
411 grey and white matter and cerebro-spinal fluid (**Ashburner and Friston, 2005**). The parameters for  
412 the affine registration of the anatomical to the MNI template were estimated, and used to normalize  
413 the grey matter probability maps of each individual to the MNI space. A group MNI source-projection  
414 grid with a resolution of 3 mm was prepared including only voxels associated with a group-average  
415 grey-matter probability of at least 0.25. The projection grid excluded various subcortical structures,  
416 identified using the AAL atlas (e.g., vermis, caudate, putamen and the cerebellum). Leadfields were  
417 computed based on a single shell conductor model. Time-domain projections were obtained on a  
418 block-by-block basis for each frequency band using LCMV spatial filters (regularization = 5%) along  
419 the dipole orientation of maximum variance.

### 420 **Analysis of speech to brain entrainment**

421 Motivated by previous work (**Gross et al., 2013; Ng et al., 2013**), we considered eight partly overlapping  
422 frequency bands (0.25–1 Hz, 1–4 Hz, 4–8 Hz, 8–12 Hz, 12–18 Hz, 18–24 Hz, 24–36 Hz, and 30–48 Hz), and  
423 isolated them from the full-spectrum MEG and speech envelope signals using a forward-reverse  
424 4<sup>th</sup> order Butterworth filter (magnitude of frequency response at band limits = -6 dB). Entrainment  
425 was quantified using the mutual information (MI) between the filtered MEG and speech-envelope  
426 time courses (**Cogan and Poeppel, 2011; Gross et al., 2013; Kayser et al., 2015b; Keitel et al., 2017; Ng  
427 et al., 2012**). The MI was calculated using a recently developed bin-less approach based on statistical  
428 copulas, which provides greater sensitivity than methods based on binned signals (**Ince et al., 2016a**).

429 To quantify the entrainment of brain activity to the speech envelope we first determined the  
430 optimal time lag between MEG signals and the acoustic stimulus for individual bands and source  
431 voxels using a permutation-based RFX estimate. Lag estimates were obtained based on a quadratic  
432 fit, excluding lags with insignificant MI (permutation-based FDR = 0.01). Voxels without an estimate  
433 were assigned the median estimate within the same frequency band, and volumetric maps of the  
434 optimal lags were smoothed with a Gaussian (FWHM = 10 mm). Speech MI was then estimated for  
435 each band and voxel using the optimal lag. The significance of group-level speech MI assessed within  
436 a permutation-based RFX framework that relied on MI values corrected for bias at the single-subject  
437 level, and on cluster mass enhancement of the test statistics corrected for multiple comparisons at  
438 the second level (**Maris and Oostenveld, 2007**). At the single-subject level, null distributions were  
439 obtained by shuffling the assignment of speech and MEG, independently for each participant, i.e.  
440 by permuting the 6 speech segments within each of the 8 experimental conditions (using the same  
441 permutation across bands). Participant-specific bias-corrected speech MI values were then defined  
442 as the actual MI minus the median MI across all 720 possible null permutations. Group-level RFX  
443 testing relied on T-statistics for the null-hypothesis that the participant-averaged bias-corrected MI  
444 was significantly larger than zero. To this end we generated 10,000 samples of the group-averaged  
445 MI from the participant-specific null distributions, used cluster-mass enhancement across voxels and  
446 frequencies (cluster-forming threshold  $T(18) = 2.1$ ) to extract the maximum cluster T across frequency  
447 bands and voxels, and considered as significant a cluster-enhanced T statistic higher than the 95<sup>th</sup>  
448 percentile of the permutation distribution (corresponding to FWE = 0.05).

449 To determine whether speech entrainment was modulated by the experimental factors we used  
450 a permutation-based RFX GLM framework (**Winkler et al., 2014**). For each participant individually  
451 we considered the condition-specific bias-corrected MI averaged across repetitions and estimated  
452 the coefficients of a GLM for predicting MI based on SNR (2, 4, 6, 8 dB), VIVN (1 = Visual Informative; -1  
453 = Visual Not informative), and their interaction. We computed a group-level T-statistic for assessing  
454 the hypothesis that the across-participant average GLM coefficient was significantly different than

455 zero, using cluster-mass enhancement across voxels and frequencies. Permutation testing relied on  
456 the Freedman-Lane procedure (**Freedman and Lane, 1983**). Independently for each participant and  
457 GLM effect, we estimated the parameters of a reduced GLM that includes all of the effects but the  
458 one to be tested and extracted the residuals of the prediction. We then permuted the condition-  
459 specific residuals and extracted the GLM coefficient for the effect of interest estimated for these  
460 reshuffled residuals. We obtained a permutation T statistic for the group-average GLM coefficient  
461 of interest using the max-statistics. We considered as significant T values whose absolute value was  
462 higher than the 95<sup>th</sup> percentile of the absolute value of 10,000 permutation samples, correcting for  
463 multiple comparisons across voxels / bands (FWE = 0.05). We only considered significant GLM effects  
464 in conjunction with a significant condition-average entrainment.

### 465 **Analysis of directed functional connectivity**

466 To quantify directed functional connectivity we relied on the concept of Wiener-Granger causality  
467 and its information theoretic implementation known as Transfer Entropy or directed information  
468 (DI, **Massey, 1990; Schreiber, 2000; Vicente et al., 2011; Wibral et al., 2011**). Directed information in  
469 its original formulation (**Massey, 1990**, termed DI\* here) quantifies causal connectivity by measuring  
470 the degree to which the past of a seed predicts the future of a target signal, conditional on the past  
471 of the target, defined at a specific lag ( $\tau_{Brain}$ ):

$$DI^* (\tau_{Brain}) = I (Target_t; Seed_{t-\tau} | Target_{t-\tau}) \quad (1)$$

472 While DI\* provides a measure of the overall directed influence from seed to target, it can be sus-  
473 ceptible to statistical biases arising from limited sampling, common inputs or signal auto-correlations  
474 (**Besserve et al., 2015, 2010; Ince et al., 2016a; Panzeri et al., 2007**). We regularized and made this  
475 measure more conservative by subtracting out values of DI computed at fixed values of speech enve-  
476 lope. This subtraction removes terms -- including the statistical biases described above -- that can-  
477 not possibly carry speech information (because they are computed at fixed speech envelope). This  
478 results in an estimate that is statistically more robust, more conservative and more directly related  
479 to changes in the sensory input than classical transfer entropy (termed directed feature information  
480 in **Ince et al., 2015, 2016a**). Practically, DI was defined here as

$$DI^* (\tau_{Brain}, \tau_{Speech}) = DI^* (\tau_{Brain}) - DI^* (\tau_{Brain}) | Speech (\tau_{Speech}) \quad (2)$$

481 where  $DI^* | Speech$  denotes the DI\* conditioned on the speech envelope. Positive values of DI indicate  
482 directed functional connectivity between seed and target at a specific brain ( $\tau_{Brain}$ ) and speech lag  
483 ( $\tau_{Speech}$ ). The actual DI values were furthermore Z-scored against random effects to further enhance  
484 the robustness of this connectivity index, which facilitates statistical comparisons between condi-  
485 tions across subjects (**Besserve et al., 2015**). To this end DI, as estimated for each participant and  
486 connection from Eq. 2, was Z-scored against the distribution of DI values obtained from condition-  
487 shuffled estimates (using the same randomization procedure as for MI). DI was computed for speech  
488 lags between 0 and 500 ms and brain lags between 0 and 250 ms, at steps of one sample (1/150 Hz).  
489 We estimated DI on the frequency range of 0.25–8 Hz (forward-reverse 4th order Butterworth filter)  
490 and by considering the bivariate MEC response defined by the band-passed source signal and its  
491 first-order difference (**Ince et al., 2016a,b**). Seeds for the DI analysis were the global and local peaks  
492 of the GLM-T maps quantifying the SNR, VIVN and SNRxVIVN modulation of entrainment, and the  
493 SFG-R voxel characterized by the peak negative effect of SNR in the visual informative condition, for  
494 a total of 8 seeds (Table 1 and Fig. 3E). To test for the significance of condition-average DI we used the  
495 same permutation-based RFX approach as for speech MI, testing the hypothesis that bias-corrected  
496  $DI > 0$ . We used 2D cluster-mass enhancement of the T statistics within speech/brain lag dimensions  
497 correcting for multiple comparisons across speech and brain lags (FWE = 0.05). To test for significant  
498 DI effects with experimental conditions we relied on the same GLM strategy as for MI effects, again  
499 with the same differences pertaining to cluster enhancement and comparison correction (FWE =  
500 0.05 across lags and seed/target pairs). We only considered DI modulations in conjunction with a  
501 significant condition-average DI.

## 502 **Neuro-behavioural correlations**

503 We used a permutation-based RFX approach to assess (1) whether an increase in condition-specific  
504 speech-MI or DI was associated with an increase in behavioural performance, and (2) whether the  
505 visual enhancement (VI-VN) of MI or DI was associated with stronger behavioural gains. We focused  
506 on the 8 regions used as seeds for the DI analysis. For speech-MI we initially tested whether the  
507 participant-average Fisher Z-transformed correlation between condition-specific performance and  
508 speech-MI was significantly larger than zero. Uncorrected p-values were computed using the per-  
509 centile method, where FWE = 0.05 p-values corrected across regions were computed using maxi-  
510 mum statistics. We subsequently tested the positive correlation between SNR-specific visual gains  
511 (VI-VN) in speech-MI and behavioural performance using the same approach, but considered only  
512 those regions characterized by a significant condition-specific MI/performance association. For DI,  
513 we focused on those lags characterized by a significant SNR, VIVN, or SNRVIVN DI modulation.  
514 Significance testing proceeded as for speech MI, except that Z-transformed correlations were com-  
515 puted independently for each lag and then averaged across lags (FWE = 0.05 corrected across all  
516 seed/target pairs).

## 517 **Acknowledgements**

518 This research was supported by the UK Biotechnology and Biological Sciences Research Council  
519 (BBSRC, BB/L027534/1). CK is supported by the European Research Council (ERC-2014-CoG; grant  
520 No 646657); BLG by BBSRC BB/M009742/1; JG by the Wellcome Trust (Joint Senior Investigator Grant,  
521 No 098433); SP by the Autonomous Province of Trento ("Grandi Progetti 2012, Characterizing and Im-  
522 proving Brain Mechanisms of Attention--ATTEND"); PGS by the Wellcome Trust (Senior Investigator  
523 Grant 107802/Z/15/Z). Competing interests: none.

## 524 **Additional information**

### 525 **Author contributions**

526 Conceptualization: CK; Methodology: BLG, RAAI, JG, SP, PGS, CK; Software: BLG, RAAI, CK; Validation:  
527 BLG, CK; Formal Analysis: BLG, CK; Investigation: BLG, CK; Resources: BLG, RAAI, CK; Data Curation:  
528 BLG; Writing -- Original Draft: BLG, CK; Writing -- Review & Editing: BLG, RAAI, JG, SP, PGS, CK; Visu-  
529 alization: BLG, CK; Supervision: CK; Project Administration: CK; Funding Acquisition: CK, JG.

## 530 **References**

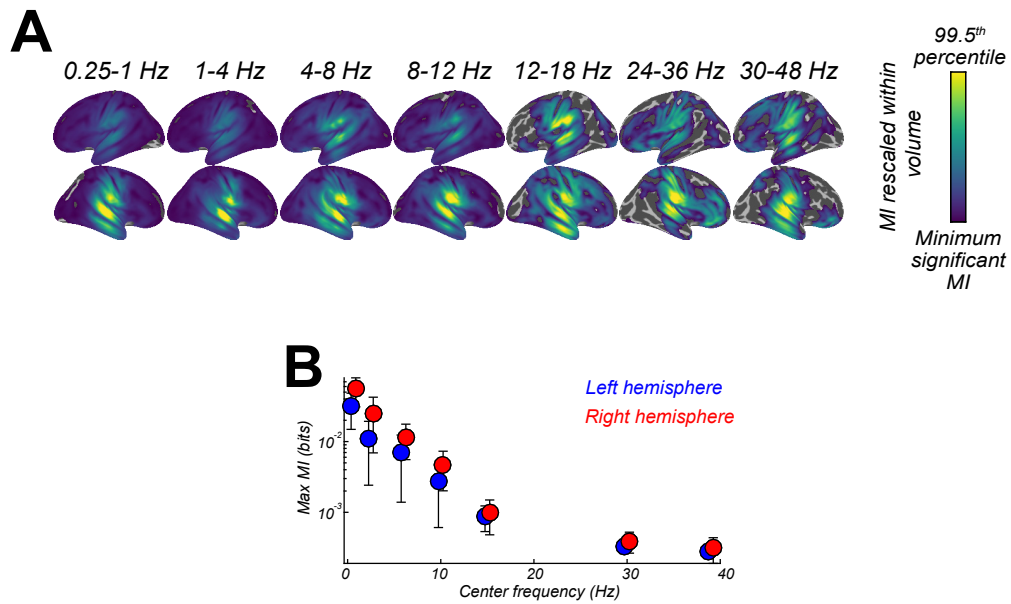
- 531 **Abrams DA**, Ryali S, Chen T, Balaban E, Levitin DJ, Menon V. Multivariate activation and connectivity patterns dis-  
532 criminate speech intelligibility in Wernicke's, Broca's, and Geschwind's areas. *Cereb Cortex*. 2013; **23**:1703-1714.
- 533 **Alho J**, Lin FH, Sato M, Tiitinen H, Sams M, Jääskeläinen IP. Enhanced neural synchrony between left auditory  
534 and premotor cortex is associated with successful phonetic categorization. *Front Psychol*. 2014; **5**:394.
- 535 **Arnal LH**, Giraud AL. Cortical oscillations and sensory predictions. *Trends Cogn Sci*. 2012; **16**:390-398.
- 536 **Arnal LH**, Morillon B, Kell CA, Giraud AL. Dual neural routing of visual facilitation in speech processing. *J Neurosci*.  
537 2009; **29**:13445-13453.
- 538 **Ashburner J**, Friston KJ. Unified segmentation. *Neuroimage*. 2005; **26**:839-851.
- 539 **Bastos AM**, Usrey WM, Adams RA, Mangun GR, Fries P, Friston KJ. Canonical microcircuits for predictive coding.  
540 *Neuron*. 2012; **76**:695-711.
- 541 **Beauchamp MS**, Argall BD, Bodurka J, Duyn JH, Martin A. Unraveling multisensory integration: Patchy organiza-  
542 tion within human STS multisensory cortex. *Nat Neurosci*. 2004; **7**:1190-1192.
- 543 **Bernstein LE**, Auer ET, Takayanagi S. Auditory speech detection in noise enhanced by lipreading. *Speech Com-*  
544 *mun*. 2004; **44**:5-18.
- 545 **Besle J**, Fischer C, Bidet-Caulet A, Lecaigard F, Bertrand O, Giard MH. Visual activation and audiovisual interac-  
546 tions in the auditory cortex during speech perception: Intracranial recordings in humans. *J Neurosci*. 2008;  
547 **28**:14301-14310.
- 548 **Besserve M**, Lowe SC, Logothetis NK, Schölkopf B, Panzeri S. Shifts of gamma phase across primary visual cortical  
549 sites reflect dynamic stimulus-modulated information transfer. *PLoS Biol*. 2015; **13**:e1002257.
- 550 **Besserve M**, Schölkopf B, Logothetis NK, Panzeri S. Causal relationships between frequency bands of extracellular  
551 signals in visual cortex revealed by an information theoretic analysis. *J Comput Neurosci*. 2010; **29**:547-566.

- 552 **Binder JR**, Desai RH, Graves WW, Conant LL. Where is the semantic system? A critical review and meta-analysis  
553 of 120 functional neuroimaging studies. *Cereb Cortex*. 2009; **19**:2767–2796.
- 554 **Bornkessel-Schlesewsky I**, Schlesewsky M, Small SL, Rauschecker JP. Neurobiological roots of language in pri-  
555 mate audition: Common computational properties. *Trends Cogn Sci*. 2015; **19**:142–150.
- 556 **Bourguignon M**, De Tiège X, Op de Beeck M, Ligot N, Paquier P, Van Bogaert P, et al. The pace of prosodic phrasing  
557 couples the listener's cortex to the reader's voice. *Hum Brain Mapp*. 2013; **34**:314–326.
- 558 **Brainard DH**. The psychophysics toolbox. *Spat Vis*. 1997; **10**:433–436.
- 559 **Callan DE**, Jones JA, Callan A. Multisensory and modality specific processing of visual speech in different regions  
560 of the premotor cortex. *Front Psychol*. 2014; **5**:389.
- 561 **Callan DE**, Jones JA, Callan AM, Akahane-Yamada R. Phonetic perceptual identification by native-and second-  
562 language speakers differentially activates brain regions involved with acoustic phonetic processing and those  
563 involved with articulatory-auditory/orosensory internal models. *Neuroimage*. 2004; **22**:1182–1194.
- 564 **Callan DE**, Jones JA, Munhall K, Callan AM, Kroos C, Vatikiotis-Bateson E. Neural processes underlying perceptual  
565 enhancement by visual speech gestures. *Neuroreport*. 2003; **14**:2213–2218.
- 566 **Chandrasekaran C**, Lemus L, Ghazanfar AA. Dynamic faces speed up the onset of auditory cortical spiking re-  
567 sponses during vocal detection. *Proc Natl Acad Sci U S A*. 2013; **110**:E4668–E4677.
- 568 **Chandrasekaran C**, Trubanova A, Stillitano S, Caplier A, Ghazanfar AA. The natural statistics of audiovisual speech.  
569 *PLoS Comput Biol*. 2009; **5**:e1000436.
- 570 **Chu YH**, Lin FH, Chou YJ, Tsai KWK, Kuo WJ, Jääskeläinen IP. Effective cerebral connectivity during silent speech  
571 reading revealed by functional magnetic resonance imaging. *PLoS One*. 2013; **8**:e80265.
- 572 **Clos M**, Langner R, Meyer M, Oechslin MS, Zilles K, Eickhoff SB. Effects of prior information on decoding degraded  
573 speech: An fMRI study. *Hum Brain Mapp*. 2014; **35**:61–74.
- 574 **Cogan GB**, Poeppel D. A mutual information analysis of neural coding of speech by low-frequency MEG phase  
575 information. *J Neurophysiol*. 2011; **106**:554–563.
- 576 **Ding N**, Simon JZ. Neural coding of continuous speech in auditory cortex during monaural and dichotic listening.  
577 *J Neurophysiol*. 2012; **107**:78–89.
- 578 **Ding N**, Simon JZ. Adaptive temporal encoding leads to a background-insensitive cortical representation of  
579 speech. *J Neurosci*. 2013; **33**:5728–5735.
- 580 **Du Y**, Buchsbaum BR, Grady CL, Alain C. Noise differentially impacts phoneme representations in the auditory  
581 and speech motor systems. *Proc Natl Acad Sci U S A*. 2014; **111**:7126–7131.
- 582 **Evans S**, Davis MH. Hierarchical organization of auditory and motor representations in speech perception: Evi-  
583 dence from searchlight similarity analysis. *Cereb Cortex*. 2015; **25**:4772–4788.
- 584 **Fernández LM**, Visser M, Ventura-Campos N, Ávila C, Soto-Faraco S. Top-down attention regulates the neural  
585 expression of audiovisual integration. *Neuroimage*. 2015; **119**:272–285.
- 586 **Ferstl EC**, Neumann J, Bogler C, Von Cramon DY. The extended language network: A meta-analysis of neuroimag-  
587 ing studies on text comprehension. *Hum Brain Mapp*. 2008; **29**:581–593.
- 588 **Fonteneau E**, Bozic M, Marslen-Wilson WD. Brain network connectivity during language comprehension: Inter-  
589 acting linguistic and perceptual subsystems. *Cereb Cortex*. 2015; **25**:3962–3976.
- 590 **Freedman D**, Lane D. A nonstochastic interpretation of reported significance levels. *J Bus Econ Stat*. 1983;  
591 **1**:292–298.
- 592 **Ghazanfar AA**, Maier JX, Hoffman KL, Logothetis NK. Multisensory integration of dynamic faces and voices in  
593 rhesus monkey auditory cortex. *J Neurosci*. 2005; **25**:5004–5012.
- 594 **Ghazanfar AA**, Schroeder CE. Is neocortex essentially multisensory? *Trends Cogn Sci*. 2006; **10**:278–285.
- 595 **Giraud AL**, Poeppel D. Cortical oscillations and speech processing: Emerging computational principles and  
596 operations. *Nat Neurosci*. 2012; **15**:511–517.
- 597 **Gow DW**, Olson BB. Sentential influences on acoustic-phonetic processing: A Granger causality analysis of mul-  
598 timodal imaging data. *Lang Cogn Neurosci*. 2016; **31**:841–855.
- 599 **Grant KW**, Seitz PF. Measures of auditory-visual integration in nonsense syllables and sentences. *J Acoust Soc*  
600 *Am*. 1998; **104**:2438–2450.
- 601 **Greenberg S**, Carvey H, Hitchcock L, Chang S. Temporal properties of spontaneous speech—a syllable-centric  
602 perspective. *Journal of Phonetics*. 2003; **31**:465–485.
- 603 **Gross J**, Hoogenboom N, Thut G, Schyns P, Panzeri S, Belin P, et al. Speech rhythms and multiplexed oscillatory  
604 sensory coding in the human brain. *PLoS Biol*. 2013; **11**:e1001752.
- 605 **Guediche S**, Blumstein S, Fiez J, Holt LL. Speech perception under adverse conditions: Insights from behavioral,  
606 computational, and neuroscience research. *Front Syst Neurosci*. 2013; **7**:126.
- 607 **Hasson U**, Ghazanfar AA, Galantucci B, Garrod S, Keysers C. Brain-to-brain coupling: A mechanism for creating  
608 and sharing a social world. *Trends Cogn Sci*. 2012; **16**:114–121.
- 609 **Hasson U**, Nusbaum HC, Small SL. Brain networks subserving the extraction of sentence information and its  
610 encoding to memory. *Cereb Cortex*. 2007; **17**:2899–2913.
- 611 **Hasson U**, Skipper JI, Nusbaum HC, Small SL. Abstract coding of audiovisual speech: Beyond sensory represen-  
612 tation. *Neuron*. 2007; **56**:1116–1126.

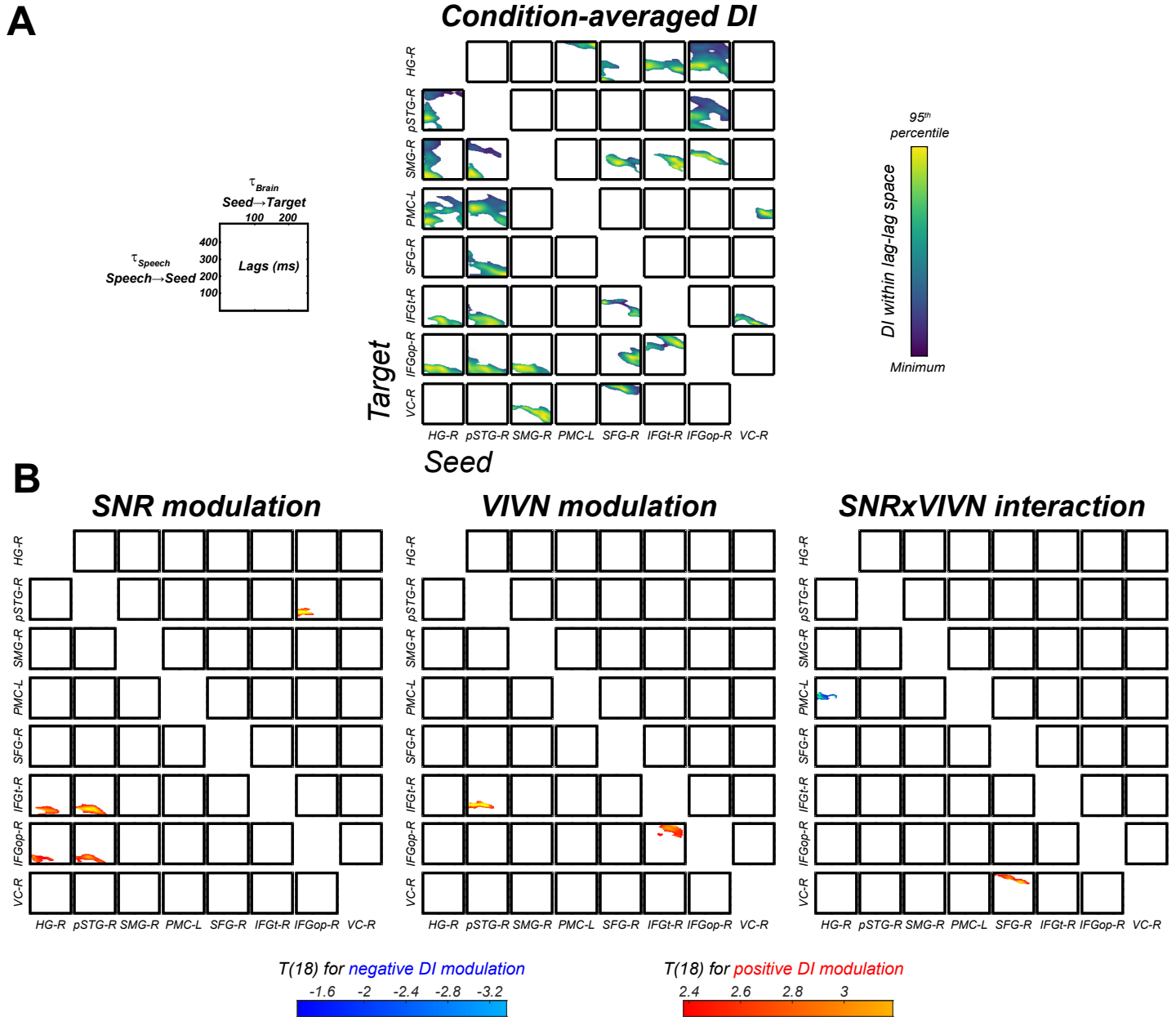


- 613 **Heim S**, Eickhoff SB, Amunts K. Specialisation in Broca's region for semantic, phonological, and syntactic fluency?  
614 *Neuroimage*. 2008; **40**:1362–1368.
- 615 **Hickok G**. Computational neuroanatomy of speech production. *Nat Rev Neurosci*. 2012; **13**:135–145.
- 616 **Hickok G**, Poeppel D. The cortical organization of speech processing. *Nat Rev Neurosci*. 2007; **8**:393–402.
- 617 **Hipp JF**, Siegel M. Dissociating neuronal gamma-band activity from cranial and ocular muscle activity in EEG.  
618 *Front Hum Neurosci*. 2013; **7**:338.
- 619 **Ince RAA**, Giordano BL, Kayser C, Rousselet GA, Gross J, Schyns PG. A statistical framework for neuroimaging  
620 data analysis based on mutual information estimated via a Gaussian copula. *Hum Brain Mapp*. 2016; . <http://doi.org/10.1002/hbm.23471>.
- 621
- 622 **Ince RAA**, Jaworska K, Gross J, Panzeri S, van Rijsbergen NJ, Rousselet GA, et al. The Deceptively Simple N170  
623 Reflects Network Information Processing Mechanisms Involving Visual Feature Coding and Transfer Across  
624 Hemispheres. *Cereb Cortex*. 2016; . <http://doi.org/10.1093/cercor/bhw196>.
- 625 **Ince RAA**, van Rijsbergen NJ, Thut G, Rousselet GA, Gross J, Panzeri S, et al. Tracing the Flow of Perceptual Features  
626 in an Algorithmic Brain Network. *Sci Rep*. 2015; **5**:17681.
- 627 **Kayser C**, Logothetis NK, Panzeri S. Visual enhancement of the information representation in auditory cortex.  
628 *Curr Biol*. 2010; **20**:19–24.
- 629 **Kayser C**, Wilson C, Safaai H, Sakata S, Panzeri S. Rhythmic Auditory Cortex Activity at Multiple Timescales Shapes  
630 Stimulus-Response Gain and Background Firing. *J Neurosci*. 2015; **35**:7750–7762.
- 631 **Kayser SJ**, Ince RAA, Gross J, Kayser C. Irregular speech rate dissociates auditory cortical entrainment, evoked  
632 responses, and frontal alpha. *J Neurosci*. 2015; **35**:14691–14701.
- 633 **Kayser SJ**, McNair SW, Kayser C. Prestimulus influences on auditory perception from sensory representations  
634 and decision processes. *Proc Natl Acad Sci U S A*. 2016; **113**:4842–4847.
- 635 **Keitel A**, Ince RAA, Gross J, Kayser C. Auditory cortical delta-entrainment interacts with oscillatory power in  
636 multiple fronto-parietal networks. *Neuroimage*. 2017; **147**:32–42.
- 637 **Krieger-Redwood K**, Gaskell MG, Lindsay S, Jefferies E. The selective role of premotor cortex in speech percep-  
638 tion: A contribution to phoneme judgements but not speech comprehension. *J Cognitive Neurosci*. 2013;  
639 **25**:2179–2188.
- 640 **Kriegeskorte N**, Goebel R, Bandettini P. Information-based functional brain mapping. *Proc Natl Acad Sci U S A*.  
641 2006; **103**:3863–3868.
- 642 **Lakatos P**, O'Connell MN, Barczak A, Mills A, Javitt DC, Schroeder CE. The leading sense: Supramodal control of  
643 neurophysiological context by attention. *Neuron*. 2009; **64**:419–430.
- 644 **Lee H**, Noppeney U. Physical and perceptual factors shape the neural mechanisms that integrate audiovisual  
645 signals in speech comprehension. *J Neurosci*. 2011; **31**:11338–11350.
- 646 **Maris E**, Oostenveld R. Nonparametric statistical testing of EEG and MEG data. *J Neurosci Methods*. 2007;  
647 **164**:177–190.
- 648 **Massey J**. Causality, feedback and directed information. In: *Proc Int Symp Inf Theory Applic (ISITA-90)*; 1990. p.  
649 303–305.
- 650 **McGettigan C**, Faulkner A, Altarelli I, Obleser J, Baverstock H, Scott SK. Speech comprehension aided by multiple  
651 modalities: Behavioural and neural interactions. *Neuropsychologia*. 2012; **50**:762–776.
- 652 **Meister IG**, Wilson SM, Deblieck C, Wu AD, Iacoboni M. The essential role of premotor cortex in speech perception.  
653 *Curr Biol*. 2007; **17**:1692–1696.
- 654 **Mesgarani N**, Chang EF. Selective cortical representation of attended speaker in multi-talker speech perception.  
655 *Nature*. 2012; **485**:233–236.
- 656 **Morillon B**, Hackett TA, Kajikawa Y, Schroeder CE. Predictive motor control of sensory dynamics in auditory active  
657 sensing. *Curr Opin Neurobiol*. 2015; **31**:230–238.
- 658 **Nath AR**, Beauchamp MS. Dynamic changes in superior temporal sulcus connectivity during perception of noisy  
659 audiovisual speech. *J Neurosci*. 2011; **31**:1704–1714.
- 660 **Ng BSW**, Logothetis NK, Kayser C. EEG phase patterns reflect the selectivity of neural firing. *Cereb Cortex*. 2013;  
661 **23**:389–398.
- 662 **Ng BSW**, Schroeder T, Kayser C. A precluding but not ensuring role of entrained low-frequency oscillations for  
663 auditory perception. *J Neurosci*. 2012; **32**:12268–12276.
- 664 **Oostenveld R**, Fries P, Maris E, Schoffelen JM. FieldTrip: Open source software for advanced analysis of MEG, EEG,  
665 and invasive electrophysiological data. *Comput Intell Neurosci*. 2010; **2011**:156869.
- 666 **Osnes B**, Hugdahl K, Specht K. Effective connectivity analysis demonstrates involvement of premotor cortex  
667 during speech perception. *Neuroimage*. 2011; **54**:2437–2445.
- 668 **Panzeri S**, Senatore R, Montemurro MA, Petersen RS. Correcting for the sampling bias problem in spike train  
669 information measures. *J Neurophysiol*. 2007; **98**:1064–1072.
- 670 **Park H**, Ince RAA, Schyns PG, Thut G, Gross J. Frontal top-down signals increase coupling of auditory low-  
671 frequency oscillations to continuous speech in human listeners. *Curr Biol*. 2015; **25**:1649–1653.
- 672 **Park H**, Kayser C, Thut G, Gross J. Lip movements entrain the observers' low-frequency brain oscillations to facili-  
673 tate speech intelligibility. *eLife*. 2016; **5**:e14521.

- 674 **Peelle JE**, Davis MH. Neural oscillations carry speech rhythm through to comprehension. *Front Psychol.* 2012;  
675 **3**:320.
- 676 **Peelle JE**, Sommers MS. Prediction and constraint in audiovisual speech perception. *Cortex.* 2015; **68**:169-181.
- 677 **Pickering MJ**, Garrod S. An integrated theory of language production and comprehension. *Behav Brain Sci.* 2013;  
678 **36**:329-347.
- 679 **Poeppl D**. The analysis of speech in different temporal integration windows: Cerebral lateralization as “asym-  
680 metric sampling in time”. *Speech Commun.* 2003; **41**:245-255.
- 681 **Poeppl D**. The neuroanatomic and neurophysiological infrastructure for speech and language. *Curr Opin Neu-*  
682 *robiol.* 2014; **28**:142-149.
- 683 **Price CJ**. A review and synthesis of the first 20 years of PET and fMRI studies of heard speech, spoken language  
684 and reading. *Neuroimage.* 2012; **62**:816-847.
- 685 **Rauschecker JP**, Scott SK. Maps and streams in the auditory cortex: Nonhuman primates illuminate human  
686 speech processing. *Nat Neurosci.* 2009; **12**:718-724.
- 687 **Riedel P**, Ragert P, Schelinski S, Kiebel SJ, von Kriegstein K. Visual face-movement sensitive cortex is relevant for  
688 auditory-only speech recognition. *Cortex.* 2015; **68**:86-99.
- 689 **Rohe T**, Noppeney U. Cortical hierarchies perform Bayesian causal inference in multisensory perception. *PLoS*  
690 *Biol.* 2015; **13**:e1002073.
- 691 **Ross LA**, Saint-Amour D, Leavitt VM, Javitt DC, Foxe JJ. Do you see what I am saying? Exploring visual enhance-  
692 ment of speech comprehension in noisy environments. *Cereb Cortex.* 2007; **17**:1147-1153.
- 693 **Schepers IM**, Yoshor D, Beauchamp MS. Electrooculography reveals enhanced visual cortex responses to visual  
694 speech. *Cereb Cortex.* 2015; **25**:4103-4110.
- 695 **Schreiber T**. Measuring information transfer. *Phys Rev Lett.* 2000; **85**:461-464.
- 696 **Schroeder CE**, Lakatos P, Kajikawa Y, Partan S, Puce A. Neuronal oscillations and visual amplification of speech.  
697 *Trends Cogn Sci.* 2008; **12**:106-113.
- 698 **Schwartz JL**, Berthommier F, Savariaux C. Seeing to hear better: Evidence for early audio-visual interactions in  
699 speech identification. *Cognition.* 2004; **93**:B69-B78.
- 700 **Schwartz JL**, Savariaux C. No, there is no 150 ms lead of visual speech on auditory speech, but a range of audio-  
701 visual asynchronies varying from small audio lead to large audio lag. *PLoS Comput Biol.* 2014; **10**:e1003743.
- 702 **Skipper JI**, Goldin-Meadow S, Nusbaum HC, Small SL. Gestures orchestrate brain networks for language under-  
703 standing. *Curr Biol.* 2009; **19**:661-667.
- 704 **Sumby WH**, Pollack I. Visual contribution to speech intelligibility in noise. *J Acoust Soc Am.* 1954; **26**:212-215.
- 705 **van Atteveldt N**, Formisano E, Goebel R, Blomert L. Integration of letters and speech sounds in the human brain.  
706 *Neuron.* 2004; **43**:271-282.
- 707 **van Wassenhove V**. Speech through ears and eyes: Interfacing the senses with the supramodal brain. *Front*  
708 *Psychol.* 2013; **4**:2.
- 709 **Vander Ghinst M**, Bourguignon M, Op de Beeck M, Wens V, Marty B, Hassid S, et al. Left Superior Temporal Gyrus  
710 Is Coupled to Attended Speech in a Cocktail-Party Auditory Scene. *J Neurosci.* 2016; **36**:1596-1606.
- 711 **Vetter P**, Smith FW, Muckli L. Decoding sound and imagery content in early visual cortex. *Curr Biol.* 2014;  
712 **24**:1256-1262.
- 713 **Vicente R**, Wibral M, Lindner M, Pipa G. Transfer entropy—a model-free measure of effective connectivity for the  
714 neurosciences. *J Comput Neurosci.* 2011; **30**:45-67.
- 715 **Wibral M**, Rahm B, Rieder M, Lindner M, Vicente R, Kaiser J. Transfer entropy in magnetoencephalographic data:  
716 Quantifying information flow in cortical and cerebellar networks. *Prog Biophys Mol Biol.* 2011; **105**:80-97.
- 717 **Wild CJ**, Yusuf A, Wilson DE, Peelle JE, Davis MH, Johnsrude IS. Effortful listening: The processing of degraded  
718 speech depends critically on attention. *J Neurosci.* 2012; **32**:14010-14021.
- 719 **Wilson SM**, Saygin AP, Sereno MI, Iacoboni M. Listening to speech activates motor areas involved in speech  
720 production. *Nat Neurosci.* 2004; **7**:701-702.
- 721 **Winkler AM**, Ridgway GR, Webster MA, Smith SM, Nichols TE. Permutation inference for the general linear model.  
722 *Neuroimage.* 2014; **92**:381-397.
- 723 **Wright TM**, Pelphrey KA, Allison T, McKeown MJ, McCarthy G. Polysensory interactions along lateral temporal  
724 regions evoked by audiovisual speech. *Cereb Cortex.* 2003; **13**:1034-1043.
- 725 **Yarkoni T**, Speer NK, Zacks JM. Neural substrates of narrative comprehension and memory. *Neuroimage.* 2008;  
726 **41**:1408-1425.
- 727 **Zion Golumbic E**, Cogan BG, Schroeder CE, Poeppel D. Visual input enhances selective speech envelope tracking  
728 in auditory cortex at a “cocktail party”. *J Neurosci.* 2013; **33**:1417-1426.



**Figure S1. Entrainment of rhythmic MEG activity to the speech envelope. (A)** Projection of significant speech MI maps, which quantify the entrainment of MEG source activity to the speech envelope, onto the Freesurfer template (FWE = 0.05; proximity = 10 mm; surface-projected significant MI maps rescaled within volume from minimum significant MI to the 99.5<sup>th</sup> percentile of the surface projection). **(B)** Peak MI in the two hemispheres as a function of frequency (mean  $\pm$  SEM).



**Figure S2. Directed functional connectivity within the speech-entrained network.** (A) Significant condition-averaged directed information (DI) values between all seed-target pairs as a function of the speech ( $\tau_{Speech}$ ) and brain lags ( $\tau_{Brain}$ ). (B): Group-level statistical maps for the GLM effects on DI of acoustic signal quality (SNR), visual informativeness (VIVN) and their interaction.