ACCEPTED MANUSCRIPT • OPEN ACCESS

Effect of visual input on syllable parsing in a computational model of a neural microcircuit for speech processing

To cite this article before publication: Anirudh Kulkarni et al 2021 J. Neural Eng. in press https://doi.org/10.1088/1741-2552/ac28d3

Manuscript version: Accepted Manuscript

Accepted Manuscript is "the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an 'Accepted Manuscript' watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors"

This Accepted Manuscript is © 2021 The Author(s). Published by IOP Publishing Ltd..

As the Version of Record of this article is going to be / has been published on a gold open access basis under a CC BY 3.0 licence, this Accepted Manuscript is available for reuse under a CC BY 3.0 licence immediately.

Everyone is permitted to use all or part of the original content in this article, provided that they adhere to all the terms of the licence https://creativecommons.org/licences/by/3.0

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions may be required. All third party content is fully copyright protected and is not published on a gold open access basis under a CC BY licence, unless that is specifically stated in the figure caption in the Version of Record.

View the article online for updates and enhancements.

Page 1 of 34

1 2		
2 3 4 5	1	Effect of visual input on syllable parsing in a computational
6 7 8	2	model of a neural microcircuit for speech processing
9 10	3	
11 12	4	Anirudh Kulkarni ¹ , Mikolaj Kegler ¹ and Tobias Reichenbach ^{1,2,*}
13 14 15	5	
15 16 17	6	¹ Department of Bioengineering and Centre for Neurotechnology, Imperial College London, South
18 19	7	Kensington Campus, SW7 2AZ, London, U.K.
20 21	8	² Department Artificial Intelligence in Biomedical Engineering, Friedrich-Alexander-Universität
22 23	9	Erlangen-Nürnberg, Konrad-Zuse-Strasse 3/5, 91056 Erlangen
24 25	10	
26 27 28	11	
29 30	12	*To whom correspondence should be addressed (email: tobias.j.reichenbach@fau.de)
31 32	13	
33 34	14	
35 36 27		
37 38 39		
40 41		
42 43		
44 45		
46 47		
48 49 50		
50 51 52		
53 54		
55 56		
57 58		
59 60	7	

15 Abstract

Seeing a person talking can help to understand them, in particular in a noisy environment. However, how the brain integrates the visual information with the auditory signal to enhance speech comprehension remains poorly understood. Here we address this question in a computational model of a cortical microcircuit for speech processing. The model consists of an excitatory and an inhibitory neural population that together create oscillations in the theta frequency range. When simulated with speech, the theta rhythm becomes entrained to the onsets of syllables, such that the onsets can be inferred from the network activity. We investigate how well the obtained syllable parsing performs when different types of visual stimuli are added. In particular, we consider currents related to the rate of syllables as well as currents related to the mouth-opening area of the talking faces. We find that currents that target the excitatory neuronal population can influence speech comprehension, both boosting it or impeding it, depending on the temporal delay and on whether the currents are excitatory or inhibitory. In contrast, currents that act on the inhibitory neurons do not impact speech comprehension significantly. Our results suggest neural mechanisms for the integration of visual information with the acoustic information in speech and make experimentally-testable predictions.

1. Introduction

Speech comprehension can benefit from other sensory input, in addition to the auditory signal, through multisensory integration [1,2]. As a striking example, seeing a speaker's face and their moving lips can improve the comprehension of speech in noise by more than 10 dB in the signal-to-noise ratio [3,4]. Such audiovisual enhancement of speech comprehension may result from different visual features such as facial gestures, hand movements, jaw movements as well as the alternating configuration of the lips, teeth, tongue, head and eyebrows [5–7]. In particular, the area of mouth opening is strongly correlated to the amplitude fluctuations in speech, and mouth movements typically precede the corresponding voice onset by about 100-300 ms [8].

Classic theories of such multisensory processing posit that primary sensory regions process only unisensory inputs [9,10]. The individual streams of information are then relayed to higher-level association cortices where the information from the various unisensory regions converge to create a multisensory percept. However, recent studies in several species have shown that the integration of auditory information with other sensory modalities can occur in the brain as early as the primary and secondary auditory cortices which were hitherto considered to be unisensory areas [11]. For instance, in adult rhesus monkeys, visual stimuli were found to modulate the activity of single neurons as well as the local field potential (LFP) in the primary auditory cortex [12–15]. Similarly, single-unit recordings as well as the LFP in the auditory cortex of anaesthetized ferrets were influenced by visual stimuli [16,17]. In awake mice, using multisite probes to sample single units across multiple cortical layers, it was demonstrated that visual stimuli influenced firing in the primary auditory cortex [18] and short-term visual deprivations led to enhanced neuronal responses and frequency selectivity to sounds in layer 4 of primary auditory cortex (A1) [19]. Experiments using a voltage-sensitive dye and optical imaging in guinea pigs observed inhibitory responses in auditory areas about 110 ms after the onset of a visual stimuli [20].

In support of multisensory processing in early sensory areas, it has further been demonstrated that direct
projections from visual areas to the auditory cortex exist in monkeys [21,22], ferrets [23], Mongolian
gerbils [24], marmosets [25] and they have also been suggested in rats [26].

In humans, using fMRI, it has similarly been found that visual stimulation and the reading of text by themselves activated the auditory cortex [27,28]. Further studies using Magnetoencephalography (MEG) showed that viewing a speaker's face improved the tracking of speech rhythms in the auditory cortex [29,30], and additionally, using intercranial electroencephalography (iEEG), it was shown that the phase of the slow oscillations in the auditory cortex could track the rhythms in a talking face [31]. Moreover, intracranial stereotactic electroencephalographic (sEEG) recordings in human patients suggest direct pathways linking early visual and auditory regions and that visual input is processed in the auditory cortex about 100 ms after the visual onset [32].

Current theories of speech processing include a role of the cortical tracking of the amplitude fluctuations in speech by the different cross-coupled neural oscillations such as delta (1 - 4 Hz), theta (4 - 8 Hz) and gamma (25 - 100 Hz) rhythms [33,34]. These oscillations occur at the rhythms set by words, syllables and phonemes, respectively. In particular, the theta band is assumed to parse speech into syllables [35– 38] thus providing temporal frames for the phonemic encoding by the gamma rhythm. A computational model of a spiking neural network for speech processing that included theta oscillations coupled to gamma oscillations showed that phonemes could indeed by decoded from the gamma activity when it was parsed by the input from the theta oscillator.

Visual enhancement of speech comprehension may, at least in part, result from the visual stimuli affecting oscillatory activity in the auditory cortex. Studies on ferrets showed indeed that information from the visual cortex was conveyed to the auditory cortex through influencing the phase of the LFP [17]. One study found increases in alpha power in the auditory cortices due to visual signals [39] whereas others observed changes, including phase resets, in the delta (3–4 Hz), theta (4–8 Hz), beta (12–30 Hz) and alpha (8–14 Hz) frequency bands [32,33,40–42].

The precise mechanisms by which visual signals can influence cortical oscillations related to speech, and thereby impact speech comprehension remain, however, elusive. Even though there have been computational models of phase resets of delta oscillations [43] and biophysical models of phase locking of oscillators [44,45], none of them investigated how these effects relate to speech processing.

In this study we employ a recently suggested model of a spiking neural network for speech processing to investigate the effect of visual input [46]. In particular, the artificial neural network includes a module for theta oscillations that can parse speech into distinct syllables. We investigate how different types of speech-related visual input influence the accuracy of the syllable parsing.

90 2. Methods

91 2.1 Architecture of the computational model

92 Our artificial spiking neural network for speech processing is based on a recently-introduced model that 93 contains coupled theta- and gamma-oscillations [46]. The theta oscillations thereby segment a speech 94 stream into individual syllables, and the neural activity in the gamma range can allow to decode the 95 syllable identity.

The auditory speech input is firstly processed by a model of the thalamus before reaching a module that produces oscillations in the theta range (figure 1a). Because we are interested in investigating the influence of slow visual input, such as related to the opening and closing of the mouth, on speech processing, our model includes only the theta oscillator and not also a gamma oscillator. When stimulated by a speech input, the spiking activity of the theta module becomes aligned to the syllable boundaries. An example speech input, its time frequency spectrogram and the resulting LFP and spiking neural activity are shown in figure 1(b).

103 The theta module produces oscillations through an interplay of excitatory neurons (Te) and inhibitory104 neurons (Ti). that are reciprocally coupled via inhibitory and excitatory synapses. The theta-band

oscillations are generated by the principle of slower feedback inhibition following fast recurrentexcitations. At the beginning of each oscillatory cycle, the excitatory input increases, resulting in an



Figure 1. Architecture of the spiking neural network and the extraction of syllable onsets. (a) Network architecture. The auditory input is decomposed through 32 frequency channels and the resulting signal is relayed through a population of relay neurons, which act as a spectro-temporal filter, to the theta module. The theta module consists of 10 excitatory neurons (Te) and 10 inhibitory neurons (Ti) and generates self-sustained oscillations in the theta frequency band. The visual input is added to either Te or Ti. (b) The theta LFP generated by an example sentence 'She had your dark suit in greasy wash water all year 'together with the estimated syllable onsets (blue, top panel). The spiking of the inhibitory Ti neurons (blue dots) is aligned to the syllable onsets (red lines, bottom

108 increase in the firing rate of the excitatory population. The inhibitory population eventually catches up

109 and brings down the firing rate of the excitatory population. As the excitatory population activity goes

110 down, and as a result the inhibitory population activity decreases, the network recovers from inhibition

and the excitatory firing rate increases again. This results in a rhythmic behaviour that is referred to as

112 Pyramidal Interneuron Theta (PIN-TH) mechanism, in analogy with the Pyramidal Interneuron Gamma

113 (PING) model [47].

114 We consider 10 excitatory neurons that are reciprocally connected to each other. Likewise, we model 115 10 inhibitory neurons with all-to-all connections as well. The all-to-all connectivity within the Te 116 neurons respectively the Ti neurons means that we model a local cortical network.

117 The neurons are modelled as leaky integrate-and-fire neurons with the following dynamics for the 118 voltage V_i for cell *i*:

$$C\frac{dV_i}{dt} = g_L(V_L - V_i) + I_i^{SYN}(t) + I_i^{Inp,aud}(t) + I_i^{Inp,vis}(t) + I_i^{DC} + \eta(t)$$
(1)

120 where *C* is the capacitance of the cellular membrane;
$$g_L$$
 and V_L are the conductance and the reversal
121 potential of the leak current; $t_1^{SYN}(t)$, $l_1^{Inpastd}(t)$, $l_1^{Inpastd}(t)$, l_t^{BC} are the synaptic current, the auditory
122 stimulus-induced current, the visual stimulus-induced current and the constant direct current delivered
123 to the cell, $\eta(t)$ is white Gaussian noise with a variance of σ^2 . Whenever the membrane potential of
124 the neuron reaches the threshold potential V_{THR} , a spike is generated and returned to the recet potential
125 V_{RESET} .
126 The synaptic current $l_{ij}^{SYN}(t)$ to the postsynaptic neuron *i* from the presynaptic neuron *j* is modelled as
127 follows:
128 $l_{ij}^{SYN}(t) = g_{Ij}s_{Ij}(t)(V_{Ij}^{SYN} - V_{I}(t))$ (2)
129 where g_{Ij} is the conductance of the synapse connecting neuron *j* to neuron *i*; $s_{Ij}(t)$ is the activation
130 variable of the synapse, and V_j^{SYN} is the equilibrium potential of the synaptic current from neuron *j*.
131 The dynamics of the activation variables $s_{Ij}(t)$ of the neurons are described by the following set of
132 equations:
133 $\frac{d_{ij}t_{II}}{dt} = -\frac{s_{Ij}}{s_{Ij}} + \delta(t - t_{Ij}^{SPN})$. (3)
134 $\frac{d_{ij}t_{II}}{dt} = \frac{s_{Ij}^{S-s_{IJ}}}{s_{Ij}^{T}}$. (4)
135 where x_{Ij}^{R} are activation variables of the synapse from neuron *j* to neuron *i*; $\delta(t - t_{J}^{SPR})$ denotes a spike
136 generation in a presymptic neuron *j* at the time t_{J}^{SYN} , and τ_{I}^{R} and τ_{I}^{R} are time constants of synaptic rise
137 and decay of the presymptic neuron *j*, respectively.
138 Therefore, $t_{Ij}^{SYN}(t)$, the sum of all synaptic inputs from the cells projecting to the *t*h neuron, is given
139 by
140 $l_{I}^{CYN}(t) = \sum_{I} g_{I_{I}}(s_{I_{I}}(t) (V_{I}^{SYN} - V_{I}(t))$. (5)

141 The local field potential at time *t*, LFP(*t*), is obtained by summing the absolute values of all the synaptic142 currents delivered to all the theta excitatory cells in the network [48].

143 The model parameters were adapted from [49]. The complete list of model parameters and their values 144 are presented in Table 1. All numerical simulations of the model were performed in a custom written 145 Python script using the packages SciPy[50] and Brian2, a Python package for implementing simulations 146 of networks of neurons [51]. We used a time step of 0.01ms in all our simulations.

147 2.2 Auditory stimuli and their processing in the model

Spoken English sentences from either the TIMIT dataset [52] or from the GRID corpus [53] were provided as the auditory input to the network. The TIMIT corpus reflects realistic listening scenarios by incorporating speakers of different accents and speech production rates. It comprises of over 6,300 phonetically-labelled sentences. The GRID corpus, on the other hand, contains both the audio and visual recordings of 34 speakers speaking 1,000 sentences each.

The material from the TIMIT corpus was used for the simulations with a pulse input current, whereas the data from the GRID corpus and the corresponding videos were used for simulations where the visual current corresponded to either the area of the mouth opening or its velocity. A silent period sampled from a uniform distribution in the range 250 to 750 ms was added to each sentence to provide variability in the onset of the sentence with respect to the intrinsic firing of the theta module. This was done in order to avoid any spurious phase-locking of the network rhythm to the speech input rate.

For each simulation, a random subset of 100 sentences were chosen as the speech input in the model simulation. Speech-shaped noise was then added to each of these speech inputs. To produce the speech shaped noise, another randomly selected sentence was picked from the TIMIT database. From the linear prediction coefficients of this second sentence, a linear filter was computed. The linear filter was then convolved with a white-noise Gaussian signal to yield the speech shaped noise. The speech signal and the resulting speech shaped noise signal were mixed at a SNR of 0 dB to produce the auditory input to the model.

The auditory input was firstly processed through the model of peripheral and subcortical auditory processing [54]. The subcortical model represented the cochlear filter bank and decomposed the input signal into 128 auditory channels with center frequencies that are logarithmically spaced between 100 and 4000 Hz [54]. A series of non-linear operations representing the neural processing in the auditory nerve and subcortical nuclei were then performed on this decomposed signal. The model was implemented in a custom written Python script based on the original MATLAB implementation [46].

The number of auditory channels was then reduced to 32 by taking every fourth channel from the 128 channels. In order to reflect the experimental observation of the entraining of endogenous theta activity in the auditory cortex to the syllabic rhythm of natural speech stimuli, the theta module was designed to generate bursts of spikes aligned to the syllabic onsets in the presented sentence. For this purpose, the 32 obtained auditory channels were convolved with a spectro-temporal filter and projected to the Te neurons. This spectro-temporal filter represented a population of relay neurons with weights that corresponded to the synaptic strengths [55]. It projected the inputs with a delay of up to 50 ms and predicted syllabic onsets (binary events) based on the data from the 32 auditory channels from up to 50 ms preceding time t, in steps of every 10 ms:

$$\hat{Y}(t) = \sum_{c=1}^{32} \sum_{\tau=-50}^{0} B(c,\tau) X(c,t+\tau).$$
(6)

Y(t) is a binary variable indicating the syllabic onsets in a sentence; $\hat{Y}(t)$ is an estimate of that variable; 183 c is the index of the auditory channel; τ is the latency in ms with respect to time t; B is a matrix of filter 184 coefficients and X is the input from auditory channel c at time t. The binary vector Y(t) was determined 185 such that it had a value of 1 at the onset of each syllable but was 0 elsewhere."

To obtain the coefficients *B* of this spectro-temporal filter, 1,000 sentences that were not subsequently used for any simulations of the network, were randomly chosen from the TIMIT corpus. These sentences were appended with a silence of 500-1,000 ms at the beginning and were processed through the above-described auditory periphery model and then downsampled to 100 Hz and concatenated to give *X*. The binary vector with the corresponding syllabic onsets were processed accordingly to obtain

Y. The coefficients *B* were then obtained by providing an optimal mapping between *X* and *Y* using sparse 192 bilinear regression [56]. Once the filter coefficients were obtained, we convolved the optimized kernel 193 with the 32 auditory channels and scaled it down, to regulate the input current, by a factor of 4.5 to 194 obtain the input to the Te neurons as in the original model by [46].

2.3 Syllable parsing in the model

196 The speech input was added in the model through a current, $I_i^{Inp,aud}$, to the excitatory neurons (Te), as 197 specified above. The visual input, on the other hand, was added to the model through the visual current 198 term $I_i^{Inp,vis}$ either to the pyramidal neurons (Te) or to the inhibitory neurons (Ti),

In the absence of any speech input, the model exhibited self-sustained theta oscillations. When auditory input was added, this signal was chunked into distinct units by the theta rhythm. In particular, these chunks were delineated by the rhythmic spike bursts in the theta inhibitory module and were considered to represent individual syllables. A theta spike burst was thereby considered to be represented by the spiking of at least two inhibitory neurons in the theta module within a time window of 15 ms. The timing of such a spike burst was then considered to be the time of the maximal firing rate of Ti neurons.

2.4 Analysis of syllable parsing in the model

To quantify the accuracy of the model's syllable parsing, we computed a distance measure between the syllable boundaries inferred from the network activity and the actual boundaries, called the parsing score. The parsing score was obtained in three steps: 1) we computed a distance metric between the model's predictions and actual syllabic onsets, 2) we subtracted a control distance from this measure and 3) we divided the net result by the number of syllables. A parsing score of 1 therefore corresponded to perfect parsing by the model and a parsing score of 0 is what one would expect by chance.

To compute the distance metric in the first step, we used the normalized Victor-Purpura spike distance
metric (VPd) [57] to quantify the overall misalignment of the predicted and the actual syllable onsets.
Misalignment can result from missed syllable onsets, misaligned onsets, or additional onsets inferred

from the network activity. The VPd is particularly suitable for this task (and commonly used in spike train analysis) because it captures all three types of misalignments. The VPd between two series of binary events is calculated as the minimum cost of transforming one series into the other using one of the three operations: insertion of an event, deletion of an event and shifting of an event. A cost parameter of 50 ms was used. Hence, when the timing difference between the predicted and the actual syllable boundaries was no more than 50 ms, the two were said to be matched. A value corresponding to the ratio of the time difference to the cost was added to the distance parameter. When they were more than 50 ms apart, the score was augmented by 1. This was then subtracted from a control score defined by the normalized VPd score between the syllabic onsets and uniformly distributed bursts of spikes in the same interval. The onset in the case of the control score calculation was chosen in the same way as the random onset of the sentence. The difference in the earlier distance score and the control score was then divided by the number of syllables of the sentences, to normalize the score across different acoustic speech inputs.

The parsing scores were obtained in the same way for every simulation irrespective of the external audio and visual current inputs. The analysis was implemented in a custom written Python script using methods from SciPy package. The significance of the parsing scores obtained in each visual input condition with respect to the no-visual condition was computed using the Wilcoxon signed-rank test [58]. We then applied the Benjamini-Hochberg correction to the obtained p-values to check for false discoveries from multiple comparisons [59]. The significance threshold for the hypothesis testing was set to p=0.05.

2.5 Extraction of mouth area from the videos

To extract visual information from the videos of the GRID corpus, and in particular the mouth area, we used a custom-written Python script. The videos of the GRID corpus typically had the face of a speaker on a blue background while the speaker recited the sentence. The videos had a frame rate of 25 Hz and the speakers' face had been aligned across all the frames of the video. Each image was cropped to a small region around the mouth. The corresponding cropping region was manually determined for each speaker and stayed the same throughout the video. The pixels of the lips were then extracted using the property that the intensity of the red hue of these pixels was generally greater than the intensities of the blue or green pixels. The image was then blurred with a Gaussian filter to remove small, isolated pixels. By extracting connected objects greater than a certain threshold, we could thus extract the outer boundary of the lips. An example of the extracted lip contour whose outer boundary is highlighted by green dots is shown in figure 3(a). The number of pixels enclosed within this outer boundary was then computed for each image to obtain the area of the open mouth. We z-scored the number of pixels and upsampled the resulting signal to the same frequency as the auditory signal used in the simulation of the model corresponding to the time step of the simulation of 0.01 ms. i.e. 100 KHz, to obtain the mouth-opening area. An example is shown in figure 3(b). To explore the influence of the magnitude of the visual input on the audiovisual speech processing, in certain simulations, we multiplied the resulting signal with a factor that we called the amplitude of the area of the mouth opening. Effectively, this corresponded to scaling the standard deviation of the mouth area signal. An amplitude of 1 was used for the current unless mentioned otherwise.

255 2.6 Extraction of velocity of the mouth-opening area from the videos

To obtain the velocity of the mouth-opening area, we computed the time difference of the number of pixels within the lip contour that we obtained for a given speaker in a video. This difference signal was then z-scored to obtain the velocity of the mouth-opening area. In certain simulations, we multiplied this resulting signal with a certain factor, that we termed the amplitude of the velocity of the mouthopening area. Effectively, this corresponded to scaling the standard deviation of the signal. An amplitude of 1 was used for the current unless otherwise mentioned. An example signal of the velocity of the mouth area obtained is shown in figure 4(a).

263 2.7 Adding visual input to this network

⁵⁵
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁷
⁵⁸
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁰
⁵⁰
⁵¹
⁵²
⁵³
⁵⁴
⁵⁵
⁵⁶
⁵⁶
⁵⁷
⁵⁸
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁰
⁵¹
⁵¹
⁵²
⁵³
⁵⁴
⁵⁵
⁵⁶
⁵⁶
⁵⁷
⁵⁷
⁵⁸
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁹
⁵⁰
⁵¹
⁵¹
⁵²
⁵³
⁵⁴
⁵⁵
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁷
⁵⁸
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁷
⁵⁸
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁷
⁵⁸
⁵⁶
⁵⁶
⁵⁷
⁵⁸
⁵⁶
⁵⁷
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁸
⁵⁶
⁵⁶
⁵⁷
⁵⁸
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁶
⁵⁷
⁵⁸
⁵⁶
⁵⁶</l

speech signal. Second, we considered a current that varied in proportion to the mouth-opening area.
This current represented an important feature of the visual stimuli. Third, we investigated a current
that was proportional to the velocity of the mouth-opening area. This current was chosen since the
visual cortex can extract motion aspects from videos.
We did not consider further, more complex spatiotemporal filters for the video signal. Unlike the
spectrotemporal auditory filter, such visual filters would most likely perform poorly, due to the much

areas of the brain remains poorly understood, we have chosen this simplified model of 'lip detection

higher dimensionality of the visual input. As the flow of visual information to the speech processing

275 rather than a more intricate representation of the visual signal.

These visual currents were added as $I_i^{Inp,vis}(t)$ to either the excitatory population or the inhibitory population of the theta network module. Moreover, the visual input current was offset in time with respect to the corresponding auditory signal such that we could investigate the effects of the different time-lagged offsets in visual current on the syllable parsing scores.

For each of the three different types of currents, we studied four conditions: a) adding an excitatory visual input current to the excitatory neurons of the theta module, b) adding an excitatory visual input current to the inhibitory neurons, c) adding an inhibitory visual input current to the excitatory neurons of the theta module, and d) adding an inhibitory visual input current to the inhibitory neurons. In all cases, we compared the resulting syllable parsing scores to the condition with no visual input current.

2.8 Computing the phase of the signal

To compute the instantaneous phase of the LFP, we used the Hilbert transform. The LFP from the theta
module was firstly filtered using a third-order lowpass Butterworth filter with a cutoff frequency of 30
Hz. The Hilbert transform was subsequently applied to the resulting signal to obtain the envelope and
phase of the signal. We then determined the phase of the signal at the syllabic onsets of the sentence.
The mean phase at the syllabic onsets was computed under the different visual input conditions and was
compared to the case with no visual input.

2.9 Computing the scalogram of the signal

To compute the scalogram of the LFP, we consider the frequency band between 1 Hz and 100 Hz and performed a Morlet continuous wavelet transform with logarithmically spaced frequencies over this frequency band. This function was implemented using the Time Frequency Misfit module in the Signal module of the Obspy package in Python [60]. Once the time-frequency coefficients of the scalogram were obtained, we squared their absolute values and averaged them over time to obtain different coefficients of the average squared scalogram as a function of frequency. This quantity represents the power spectral density in the LFP [61].

3. Results

We first verified that the theta module yielded oscillations in the theta frequency range. We found that,
before the beginning of a sentence, the module produced bursts at an interval of about 150 ms (figure
1(b)). These regular bursts of neuronal spikes were also visible in the LFP.

When a speech stimulus was presented to the network, the spiking activity of the theta network became
aligned to the syllable onsets (figure 1(b)). This allowed us to investigate how well this syllable parsing
through the spikes of the theta module performed, and how this performance changed with the addition
of different visual stimuli.

308 3.1 Syllable parsing score for pulsed input current

We first considered a current that consisted of pulses of a duration of 25 ms and an amplitude of 10 pA,
unless mentioned otherwise. Each pulse occurred at the onset of a syllable, although we also considered
different time lags between the pulses and the corresponding syllables. An example of an excitatory
pulse current where the onsets of the pulses coincided with the syllable onsets is shown in figure 2(a).

The parsing score for speech without any background auditory noise was 0.08, similar to the score obtained in the original model by Hyafil et al.[46]. In the audio-only condition in our simulations, the audio is comprised of a speech input with background speech-shaped noise at an SNR of 0 dB. This

resulted in a parsing score of 0.06 for the audio-only condition. These comparatively low parsing scores
reflected frequent missed syllable onsets, misaligned onsets, and additional onsets inferred from the
neural activity.

Adding an excitatory pulse current to the excitatory neurons significantly changed the parsing score (figure 2(b)). When the onset of the pulsed coincided with the syllabic onsets, the parsing score improved significantly compared to the audio-only score of 0.06, that is, compared to the condition without any visual input current. On the other hand, delays of around -75 ms as well as around 100 ms led to significantly worse syllable parsing. A positive delay hereby meant that the current pulses occurred after the corresponding syllabic onsets.

When an inhibitory pulse input was presented to the excitatory population, we observed a significant improvement in the mean parsing score at a delay of about -125 ms, as well as a significant worsening of syllable parsing at a delay of about -25 ms (figure 2(b)).

Adding an excitatory or inhibitory pulse current to the inhibitory neurons of the theta module, however, had no significant effect on the syllable parsing (figure 2(b)). This could have resulted from the recurrent inhibitory connections in the inhibitory population of the network, that may have effectively stunted the activity of the neurons in spite of an external input current.

The parsing scores showed a periodicity of about 150 ms as a function of the delay of the input pulse current when the latter was presented to the excitatory neurons. This periodicity was comparable to the periodicity in the LFP of the theta module, as evident in the autocorrelation of the LFP without a visual or speech input (figure 2(c)). Adding an excitatory pulse current at the syllabic onset presumably made the neurons ready to fire at the syllabic onset. An inhibitory current, on the other hand, reset the excitatory population, such that the neurons were ready to fire together in the next theta cycle, at the syllabic onset.

339 Next, we investigated the effect of the amplitude and duration of the input current to the excitatory
automatical and an environment of the input current in the input current: a delay of
automatical and an environment of the input current in the input

> 341 25 ms in the case of the excitatory input, and an advance of 125 ms for the inhibitory pulses. These time 342 lags were chosen because they produced the largest significant improvements in the parsing score for 343 the respective currents. The parsing scores improved with the amplitude of the visual current, in 344 particular for smaller currents below 3 pA (figure 2(d)).

To vary the duration of the pulses, we fixed the onset of the pulse current at a delay of 25 ms for the excitatory input, and at an advance of 125 ms for the inhibitory current. We then varied the location of the offset, thus varying the duration of the pulse. The mean parsing score improved as the duration of the pulse increased until 25 ms and then reduced again for longer durations (figure 2(e)). The pulse current presumably reset the activity of the population, and the reset may have been more efficient for longer pulses. However, the longer each pulse lasted, i.e. the further the offset of the pulse current was, the more delayed the reset of the theta population likely was, thus delaying the matching of the theta prediction with respect to the actual syllabic onset. This effect may have caused the degradation of the parsing score for longer pulses.



Figure 2. Effect of a pulse input current on the parsing scores. (a) An example of an excitatory pulse input current signal with the onsets of the pulses located at the syllable onsets. (b) An excitatory input to the excitatory neurons (blue) can both improve the syllable parsing or impede it, depending on the delay. A positive delay hereby means that the pulses occur after the corresponding syllable onset. An inhibitory current to the excitatory neurons can influence the syllable parsing as well (red). Neither excitatory (green) nor inhibitory (cyan) current projected to the inhibitory neurons, however, has a significant effect on the syllable parsing. (c) The autocorrelation of the LFP in the absence of a speech stimulus or a visual current shows a periodicity of about 150 ms. (d) The mean parsing scores as a function of the amplitude of the pulse current. The excitatory inputs (blue) occur at 25 ms after the syllable onsets whereas the inhibitory input (red) is presented 125 ms before the syllable onset. (e) The mean parsing scores as a function of the pulse current. The excitatory inputs (blue) are presented 25 ms after the syllable onsets whereas the inhibitory input (red) occurs 125 ms before the syllable onset. (f) The mean parsing score is optimal for a particular phase of the LFP at the syllable onset, both for the excitatory current (blue) and for the inhibitory current (red). Statistical significance is denoted by asterisks (p < 0.05, FDR correction for multiple comparisons).

356 To explicitly test this hypothesis, we computed the mean phase of the LFP at the syllabic onset for the

⁵⁵₅₆ 357 different pulse current stimulations. We then related the phase of the LFP to the parsing score (figure

358 2(f)). The parsing scores showed a strong dependency on the phase. In particular, the parsing score

improved most when the phase of the oscillation was reset to about 270° for the excitatory input, and
to about 200° for the inhibitory current.

361 3.2 Mouth-opening area

Next, we added a visual input current that corresponded to the area of the mouth (figure 3(a),(b)). We thereby considered both positive and negative amplitudes. A positive amplitude hereby meant an excitatory current, and a negative amplitude an inhibitory current. We also considered different delays between the mouth-opening current and the speech signal.

366 Presenting both an excitatory or an inhibitory current corresponding to the mouth-opening area to the367 theta excitatory population could increase as well as decrease the parsing score, depending on the delay

368 (figure 3(c)). In particular, an excitatory current led to a worsening of syllable parsing at delays of
around 50 ms. An inhibitory current at a delay of about 50 ms led to enhanced syllable parsing, whereas
delays of around -150 ms and 150 ms led to lower parsing scores.

The observed temporal dependencies resembled the ones obtained in a computational model on neurostimulation with the speech envelope [49]. This similarity may result from the considerable correlation of the mouth-opening area and the speech envelope of 0.4 [8].

374 In contrast, adding such a current to the theta inhibitory population did not affect the parsing score.

We also investigated how the improvements in the parsing score for the current presented to the excitatory neurons varied with the amplitude of the current (figure 3(d)). We thereby considered a delay of -125 ms for the excitatory current, and a delay of 50 ms for the inhibitory one. The excitatory current only produced a significant change in the parsing score at the highest amplitude. In contrast, the inhibitory current showed significantly improved syllable parsing only for small amplitudes. The latter effect may imply that this current must be of the same order as that of the speech signal in order to improve syllable parsing.



Figure 3. Effect of a current proportional to the mouth-opening area on the parsing scores. (a) The area of mouth opening is derived from the contour of the lips (green) (b) The area is then computed for every image in a video and z-scored to yield the time-varying mouth-opening area. (c) Adding an excitatory (blue) current to the excitatory neurons could significantly worsen the parsing score, whereas an inhibitory current (red) could both improve and worsen it. In contrast, neither an excitatory nor an inhibitory current presented to the inhibitory neurons had an impact on the paring scores (green and cyan). (d) The mean parsing scores as a function of the amplitude of the visual current. The excitatory input (blue) is presented 125 ms preceding the auditory input whereas the inhibitory input (red) is added at a delay of 50 ms. Statistical significance is denoted by asterisks (p<0.05, FDR correction for multiple comparisons).

3.3 Velocity of the mouth opening area

The third set of visual stimuli that we considered corresponded to the velocity of the mouth-opening area (figure 4(a)). We multiplied this current with a certain value that we refer to as the amplitude of the current. A positive amplitude results in an excitatory current, and a negative amplitude in an inhibitory one. This current was also offset in time by different delays with respect to the corresponding auditory speech input.

As for the two other types of current, we found that both excitatory and inhibitory currents presented to
 391 the excitatory neurons could increase as well as decrease the parsing scores (figure 4(b)). In particular,
 an excitatory current at a delay of about -50 ms increased the parsing score, whereas a delay of about

393 100 ms led to a decrease. An inhibitory current could enhance syllable parsing at a delay of 125 ms and394 worsen the syllable parsing when presented at a delay of about -50 ms.

When presented to the inhibitory neurons, however, such currents had no significant effect on syllableparsing (figure 4(b)).

We also explored how the amplitudes of the currents, presented to the excitatory neurons, influenced the parsing scores (figure 4(c)). We thereby considered a delay of -50 ms for the excitatory current, and a delay of 125 ms for the inhibitory current. As we observed for the current that was based on the mouthopening area, large amplitudes of the current degraded the parsing score.

401 3.4 Firing rates of the excitatory neurons under different visual input conditions

As detailed above, we found that different types of visual currents can enhance syllable parsing when presented to the excitatory neurons. In particular, these current were 1) an excitatory pulse current with a delay of 25 ms, 2) an inhibitory pulse current at a delay 125 ms, 3) an excitatory current corresponding to the mouth-opening area at a delay of -125 ms, 4) an inhibitory current corresponding to the mouthopening area at a delay of 50 ms, 5) an excitatory current corresponding to the velocity of the mouthopening area at a delay of -25 ms, and 6) an excitatory current corresponding to the velocity of the mouth-opening area at a delay of 125 ms.

We wondered how the firing rates of the excitatory neurons changed during the presentation of these currents (figure 4(d)). We found that all currents produced changes in the firing rates as compared to the lack of a visual current. Most currents led to moderately higher firing rates. However, the excitatory pulse current caused a much larger firing rate, more than twice the one obtained without visual input The inhibitory pulse current, on the other hand, yielded a somewhat lower firing rate than without visual current.



Figure 4. Effect of a current based on the velocity of the mouth-opening area on the parsing scores. (a) An example current signal. (b) Presenting an excitatory current (blue) or an inhibitory current (red) to the excitatory neurons could, at certain delays, significantly improve the parsing score as compared to no visual input (black). However, adding either an excitatory or inhibitory current to the inhibitory neurons (cyan and green) did not influence the syllable parsing. (c) The mean parsing scores as a function of the amplitude of the visual current. The excitatory inputs (blue) were presented at a delay of -25 ms, whereas the inhibitory input (red) was added at a delay of 125 ms. (d) The firing rates of the excitatory neurons under the different conditions. The excitatory pulse inputs (Pulse E) were presented at a delay of 25 ms, whereas the inhibitory pulse inputs (Pulse E) were presented at a delay of 25 ms, whereas the inhibitory current at a delay 125 ms. The excitatory current input (Area I) occurred at a delay of 50 ms. The excitatory current (Vel I) occurred at a delay of 125 ms. All input currents cause significantly different firing rates in the excitatory neurons, and in particular for the case of the excitatory pulse inputs (Pulse E). Statistical significance is denoted by asterisks (p < 0.05, FDR correction for multiple comparisons).

45 418 3.5 Spectrogram and scalogram data for the different conditions

As another assessment of the effects of the visual currents on the network activity, we investigated the LFP as well. We computed the squared absolute values of the spectrogram and determined how they varied as a function of frequency for the different input visual currents (figure 5). This quantity is the analogue of the power spectral density for the wavelet transform. An example LFP, the corresponding time-frequency spectrogram and the average of the squared absolute value of the spectrogram in the case of no-visual input is shown in figure 5(a).

We then investigated the impact of the six different currents described above that enhanced syllable parsing: 1) the excitatory pulse current with a delay of 25 ms, 2) the inhibitory pulse current at a delay 125 ms, 3) the excitatory current corresponding to the mouth-opening area at a delay of -125 ms, 4) the inhibitory current corresponding to the mouth-opening area at a delay of 50 ms, 5) the excitatory current corresponding to the velocity of the mouth-opening area at a delay of -25 ms, and 6) the excitatory current corresponding to the velocity of the mouth-opening area at a delay of 125 ms.

We found that excitatory pulses increased the overall power of the signal and shifted the location of the
maximum from around 6 Hz to 5 Hz while adding a second local maximum at around 12 Hz (figure
5(b)). The inhibitory pulse current also increased the overall power of the signal, though to a smaller
extent than the excitatory current, and shifted the maximum slightly to a lower frequency.

The excitatory and inhibitory currents that were based on the mouth-opening area both redistributed the
power of the LFP and shifted the location of the maximum to a slightly higher frequency while causing
an additional larger peak at a lower frequency (figure 5(c)). The amplitude of the maximum at the lower

438 frequency was slightly higher for the case of the excitatory current than that of the inhibitory current.

Regarding the currents based on of the velocity of the mouth-opening area, both excitation and
inhibition increased the power of the LFP and caused an additional maximum at a low frequency (figure
5(d)).



Figure 5. The mean squared scalogram for different visual currents to the excitatory neurons. (a) The spectrogram derived from an exemplary LFP s and the corresponding squared scalogram as a function of frequency. (b) The squared scalograms for three different conditions: no visual input (green), an excitatory pulse input (blue) and inhibitory pulse inputs (red). The excitatory inputs (blue) were added at 25 ms after the syllable onsets whereas the inhibitory input (red) was presented 125 ms before the syllable onset. (c) The squared scalograms for currents based on the mouth-opening area. The excitatory input (blue) preceded the auditory signal by 125 ms, whereas the inhibitory input (red) had a delay of 50 ms. (d) The mean squared scalograms for currents based on the velocity of the mouth-opening area. The excitatory signal by 25 ms whereas the inhibitory current (red) was presented with a delay of 125 ms.

4. Discussion

We studied the effects of visual input on syllable parsing in an artificial neural network for speech processing. The neural network contained a theta module that consisted of coupled excitatory as well as inhibitory neurons and produced rhythmic bursts of spikes in the theta frequency range. When stimulated by speech, the spike bursts became aligned to the syllable boundaries, parsing the speech stream into distinct functional units.

We designed the computational model to explore possible mechanisms of audio-visual integration in speech processing and generate testable hypotheses for experimental studies. The values proposed in Table 1 are a set of 'default' parameters, which may be furthermore modified if required. These values were previously used to systematically explore speech-in-noise processing in the model and were found to be a good fit showing similar trends to psychometric curves of human speech-in-noise comprehension[49]. Due to the relatively small size and low computational complexity, the model allows to quickly screen a large space of hyperparameters (as we did here) to predict the effects of different conditions on the model behaviour.

460 We investigated how the accuracy of this syllable parsing changed when an additional current was 461 added that mimicked different aspects of an accompanying visual signal. In particular, we added three 462 different types of visual input currents to the network: a pulse current, a current corresponding to the 463 mouth-opening area of the speaker, and a current corresponding to the velocity of the mouth-opening 464 area.

We found that adding each of the three types of visual input currents could enhance as well as impede
syllable parsing. However, syllable parsing was only affected when the current acted on the excitatory,
but not on the inhibitory neurons. We suppose that this is due to the recurrent inhibitory connections in
the inhibitory population of the network which stunt the activity of the neurons in the presence of an
external input current.

In the case of the pulse current, we observed that the parsing score as a function of the time lag exhibited some periodicity with a time period of 150 ms, which corresponded roughly to the time period of the theta oscillation. Furthermore, the dependency of the parsing score on the audiovisual time delay for the inhibitory pulse current was shifted with respect to the dependency for the excitatory current by about 100 ms, which suggested that the inhibitory currents inhibited the population which recovers after a theta cycle to be ready for the syllabic onset in the next theta cycle. Furthermore, in all these cases, we found that adding the visual current significantly improved the parsing score only at certain time lags of the visual current with respect to the auditory input current. This was because adding a visual input reset the phase of the theta LFP signal. The parsing score was greater if the LFP had an optimal phase at a syllable onset than when it had a non-optimal phase.

480 Regarding the current based on the mouth-opening area, a significant improvement in the parsing score
481 resulted when the visual current was inhibitory and delayed with respect to the auditory input by 50 ms,
482 as well as when it was excitatory and had an advance of 125 ms. The current based on the velocity of
483 the mouth area current could improve the syllable parsing when it was inhibitory and delayed by 150
484 ms.

When studying time lags between the auditory and the visual signal, we need to consider two components: the physical delay and the neural delay. The physical delay is the time difference between the onset of the audio signal and the visual signal, and the neural delay is the difference between the times it takes for the audio and visual input to reach the auditory cortex. The visual stimuli typically precede the auditory stimuli by about 100-300 ms, such as for the mouth movements of a speaker compared to the actual voice onset [8]. On the other hand, IERP and voltage dye recordings indicate that the visual current arrives about 100 ms in the auditory cortex after the onset of the visual signal [32]. These two results indicate that the visual input may stimulate the auditory cortex about 100 ms before the auditory input does. Our finding that inhibitory pulse input preceding the syllabic onset by about 125-150 ms can enhance syllable parsing may therefore be particularly relevant.

We have studied both excitatory as well as inhibitory currents. However, a study in guinea pig auditory cortex using a voltage-sensitive dye showed inhibitory responses about 110 ms after the onset of the visual stimuli [20]. Another study in humans using fMRI found mainly suppressed activations in auditory cortices in response to visual stimulation [28]. Together with the likely earlier activation of the auditory cortex from visual rather than from auditory input, this suggests the enhancement of syllable parsing through an inhibitory, preceding pulse current may serve as a good model for understanding how visual inputs can enhance speech comprehension.

For the currents based on the mouth-opening area and on the corresponding velocity, the physical delay between the onset of the visual signal and the auditory signal is already incorporated in the signals. Therefore, we only need to account for the neural delay in the simulations. Considering a delay of 100 ms between the auditory and the visual signals, as suggested by experiments as described above, we find an improvement in the syllable parsing from a current based on the velocity but not from a current based on the actual mouth-opening area. The important feature might therefore be the mouth-opening area velocity rather than the mouth-opening area itself. Indeed, primary visual cortex is known to behave as an edge detector and a motion detector, responsible for computing the motion of objects across scenes [62].

When investigating which temporal lags, for a particular type of current, led to enhanced syllable parsing, we found lag regions with a width of about 50 - 100 ms. Future studies could try to design other currents with wider temporal regions for enhancing syllable parsing as seen in experimental studies [63]. Furthermore, the simulations could be made for different playback rates of audiovisual input and compared with experimental input [64].

516 Because syllable parsing in humans cannot be measured behaviourally, our computational results 517 cannot be compared directly to behavioural data. However, syllable parsing is required for syllable 518 decoding, and the latter can be tested in experiments on speech comprehension. Moreover, the model 519 predictions on neural activity could be tested in neuroimaging experiments. Experimental data could 520 compare the power spectral density of EEG waves and the firing rates to see how they correspond to

2 3 4	521 the simulations on the spectra of the LFP that we have done here to further shed light on the							
- 5 6	522	2 multisensory mechanism of audiovisual processing in the brain. Further studies could also tell us h						
7 8	523	visual speech affeo	cts the different oscillatory bands spatiotemporally across the auditory	v cortex and				
9 10	524	compare the results with the experimental data [65]. Integrating such neural data with behavioral data						
11 12	525	on speech compre-	hension in a computational model will further clarify the neural mech	anisms of				
13 14	526	audiovisual speech	n processing.					
15 16 17 18	527							
10 19 20	528	Table 1. Model p	parameters	1				
20 21 22		Parameter	Description	Value				
23 24			Neuron model					
25 26		С	Cell membrane capacitance	1 pF				
27 28		V _{THR}	Spiking threshold	-40 mV				
29 30 31		V _{RESET}	Resting potential	-87 mV				
32 33		V_E^{SYN}	Equilibrium potential of excitatory neurons	0 mV				
34 35		V _I ^{SYN}	Equilibrium potential of inhibitory neurons	-80 mV				
36 37			PINTh network					
38 39 40		$g_{\scriptscriptstyle LE}$	Leak conductance in Te neurons	0.0264 nS				
41 42		$g_{_{LI}}$	Leak conductance in <i>Ti</i> neurons	0.1 nS				
43 44 45		$ au_{Te}^R$	Synaptic rise constant of Te neurons	4 ms				
45 46 47		$ au_{Ti}^R$	Synaptic rise constant of <i>Ti</i> neurons	5 ms				
48 49		$ au_{Te}^{D}$	Synaptic decay constant of Te neurons	24.3 ms				
50 51 52		$ au_{Ti}^D$	Synaptic decay constant of Ti neurons	30.36 ms				
53 54		I ^{DC} _{Te}	Constant current delivered to Te neurons	1.25 pA				
55 56		I ^{DC} I ^{Ti}	Constant current delivered to Ti neurons	0.0851 pA				
57 58 59		σ_{Te}	Variance of the noise term in T <i>e</i> neurons	0.282 pA.ms ^{1/2}				
60		77						
		\mathbf{Y}						

σ_{Ti} $g_{Te,Ti}$ $g_{Ti,Te}$ $g_{Ti,Ti}$ References 1] Alais D, N	Variance of the noise term in $T\dot{i}$ neuronsConnectivity $Ti \rightarrow Te$ synaptic conductance strength $Te \rightarrow Ti$ synaptic conductance strength $Ti \rightarrow Ti$ synaptic conductance strength	2.028 pA.ms ^{1/2} 2.07/N _{Ti} nS 3,33/N _{Te} nS 4.32/N _{Té} nS
$g_{Te,Ti}$ $g_{Ti,Te}$ $g_{Ti,Ti}$ References 1] Alais D, N	$Ti \rightarrow Te$ synaptic conductance strength $Te \rightarrow Ti$ synaptic conductance strength $Ti \rightarrow Ti$ synaptic conductance strength	2.07/N _{Ti} nS 3,33/N _{Te} nS 4.32/N _{Ti} nS
$g_{Te,Ti}$ $g_{Ti,Te}$ $g_{Ti,Ti}$ References 1] Alais D, N	$Ti \rightarrow Te$ synaptic conductance strength $Te \rightarrow Ti$ synaptic conductance strength $Ti \rightarrow Ti$ synaptic conductance strength	2.07/N _{Ti} nS 3.33/N _{Te} nS 4.32/N _{Ti} nS
$g_{Ti,Te}$ $g_{Ti,Ti}$ References 1] Alais D, N	$Te \rightarrow Ti$ synaptic conductance strength $Ti \rightarrow Ti$ synaptic conductance strength	3,33/N _{Te} nS 4.32/N _{Té} nS
<i>g</i> _{Ti,Ti} References 1] Alais D, N	$Ti \rightarrow Ti$ synaptic conductance strength	4.32/N _{Tť} nS
References		
References		5
1] Alais D, N		
1] Alais D, N		
	Newell F N and Mamassian P 2010 Multisensory proc	essing in review: From
physiology	to behaviour vol 23	
2] Opoku-Ba	ah C, Schoenhaut A M, Vassall S G, Tovar D A, Ran	nachandran R and Wallace M
T 2021 Vis	sual Influences on Auditory Behavioral, Neural, and F	Perceptual Processes: A
Review JA	RO - J. Assoc. Res. Otolaryngol.	
3] Benoit C,	Mohamadi T and Kandel S 1994 Effects of Phonetic	Context on Audio-Visual
Intelligibil	ity of French J. Speech Hear. Res. 37 1195–203	
4] Sumby W	H and Pollack I 1954 Visual Contribution to Speech	Intelligibility in Noise J.
Acoust. So	c. Am. 26 212–5	
5] Schroeder	C E, Lakatos P, Kajikawa Y, Partan S and Puce A 20	08 Neuronal oscillations and
visual amp	lification of speech Trends Cogn. Sci. 12 106–13	
6] Campbell	R 2008 The processing of audio-visual speech: Empirication Empirication and the speech speech set of the speech set of t	rical and neural bases <i>Philos</i> .
Trans. R. S	loc. B Biol. Sci. 363 1001–10	
7] Munhall F	G, Jones J A, Callan D E, Kuratate T and Vatikiotis	-Bateson E 2004 Visual
Prosody ar	d Speech Intelligibility: Head Movement Improves A	uditory Speech Perception
Psychol. S	zi. 15 133–7	
8] Chandrase	karan C, Trubanova A, Stillittano S, Caplier A and G	hazanfar A A 2009 The
2 3 4 5 6 6 7 7	 Popoku-Ba T 2021 Vis Review JA Benoit C, 1 Intelligibili Sumby W Acoust. Soc Schroeder visual amp Campbell Trans. R. S Munhall K Prosody an Psychol. Soc Chandrase natural stat 	 Opoku-Baah C, Schoenhaut A M, Vassall S G, Tovar D A, Ran T 2021 Visual Influences on Auditory Behavioral, Neural, and F Review <i>JARO - J. Assoc. Res. Otolaryngol.</i> Benoit C, Mohamadi T and Kandel S 1994 Effects of Phonetic of Intelligibility of French <i>J. Speech Hear. Res.</i> 37 1195–203 Sumby W H and Pollack I 1954 Visual Contribution to Speech <i>Acoust. Soc. Am.</i> 26 212–5 Schroeder C E, Lakatos P, Kajikawa Y, Partan S and Puce A 20 visual amplification of speech <i>Trends Cogn. Sci.</i> 12 106–13 Campbell R 2008 The processing of audio-visual speech: Empi <i>Trans. R. Soc. B Biol. Sci.</i> 363 1001–10 Munhall K G, Jones J A, Callan D E, Kuratate T and Vatikiotis Prosody and Speech Intelligibility: Head Movement Improves <i>A</i> <i>Psychol. Sci.</i> 15 133–7 Chandrasekaran C, Trubanova A, Stillittano S, Caplier A and G natural statistics of audiovisual speech <i>PLoS Comput. Biol.</i> 5

1 2			
3 4 5	551	[9]	Felleman D J and Van Essen D C 1991 Distributed hierarchical processing in the primate
5 6	552		cerebral cortex Cereb. Cortex 1 1–47
7 8	553	[10]	Treisman A M and Gelade G 1980 A Feature-Integration Theory of Attention Cogn. Psychol.
9 10	554		12 97–136
11 12	555	[11]	Schroeder C E and Foxe J 2005 Multisensory contributions to low-level, "unisensory"
13 14 15	556		processing Curr. Opin. Neurobiol. 15 454–8
15 16 17	557	[12]	Kayser C, Petkov C I and Logothetis N K 2008 Visual modulation of neurons in auditory
17 18 19	558		cortex Cereb. Cortex 18 1560–74
20 21	559	[13]	Kayser C, Petkov C I and Logothetis N K 2009 Multisensory interactions in primate auditory
22 23	560		cortex: fMRI and electrophysiology Hear. Res. 258 80-8
24 25	561	[14]	King A J and Walker K M M 2012 Integrating information from different senses in the
26 27 28	562		auditory cortex Biol. Cybern. 106 617–25
28 29	563	[15]	Ghazanfar A A, Chandrasekaran C and Logothetis N K 2008 Interactions between the
30 31	564		superior temporal sulcus and auditory cortex mediate dynamic face/voice integration in rhesus
32 33	565		monkeys J. Neurosci. 28 4457–69
34 35 36 37 38	566	[16]	Bizley J K, Nodal F R, Bajo V M, Nelken I and King A J 2007 Physiological and anatomical
	567		evidence for multisensory interactions in auditory cortex Cereb. Cortex 17 2172-89
30 39 40	568	[17]	Atilgan H, Town S M, Wood K C, Jones G P, Maddox R K, Lee A K C and Bizley J K 2018
41 42	569		Integration of Visual Information in Auditory Cortex Promotes Auditory Scene Analysis
43 44	570		through Multisensory Binding Neuron 97 640-655.e4
45 46	571	[18]	Morrill R J and Hasenstaub A R 2018 Visual information present in infragranular layers of
47 48	572		mouse auditory cortex J. Neurosci. 38 2854–62
49 50	573	[19]	Meng X, Kao J P Y, Lee H K and Kanold P O 2017 Intracortical circuits in thalamorecipient
51 52	574		layers of auditory cortex refine after visual deprivation eNeuro 4 1-11
53 54	575	[20]	Kubota M, Sugimoto S, Hosokawa Y, Ojima H and Horikawa J 2017 Auditory-visual
55 56	576	(integration in fields of the auditory cortex Hear. Res. 346 25-33
57 58	577	[21]	Falchier A, Schroeder C E, Hackett T A, Lakatos P, Nascimento-Silva S, Ulbert I, Karmos G
60	578	7	and Smiley J F 2010 Projection from visual areas V2 and prostriata to caudal auditory cortex

2 3	579		in the monkey Cereb. Cortex 20 1529–38
4 5 6	580	[22]	Smiley J F and Falchier A 2009 Multisensory connections of monkey auditory cerebral cortex
6 7	581		Hear. Res. 258 37–46
o 9 10	582	[23]	Bizley J K and King A J 2009 Visual influences on ferret auditory cortex <i>Hear. Res.</i> 258 55-
10 11 12	583		63
13 14	584	[24]	Budinger E and Scheich H 2009 Anatomical connections suitable for the direct processing of
15 16	585		neuronal information of different modalities via the rodent primary auditory cortex Hear. Res.
17 18 10	586		258 16–27
20 21	587	[25]	Cappe C and Barone P 2005 Heteromodal connections supporting multisensory integration at
22 23	588		low levels of cortical processing in the monkey Eur. J. Neurosci. 22 2886–902
24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48	589	[26]	Stehberg J, Stehberg J, Dang P T and Frostig R D 2014 Unimodal primary sensory cortices
	590		are directly connected by long-range horizontal projections in the rat sensory cortex Front.
	591		Neuroanat. 8 1–19
	592	[27]	Pekkola J, Ojanen V, Autti T, Jääskeläinen I P, Möttönen R, Tarkiainen A and Sams M 2005
	593		Primary auditory cortex activation by visual speech: An fMRI study at 3 T Neuroreport 16
	594		125–8
	595	[28]	Gau R, Bazin P L, Trampel R, Turner R and Noppeney U 2020 Resolving multisensory and
	596		attentional influences across cortical depth in sensory cortices Elife
	597	[29]	Zion Golumbic E, Cogan G B, Schroeder C E and Poeppel D 2013 Visual input enhances
	598		selective speech envelope tracking in auditory cortex at a "Cocktail Party" J. Neurosci. 33
	599		1417–26
	600	[30]	Park H, Kayser C, Thut G and Gross J 2016 Lip movements entrain the observers' low-
49 50	601		frequency brain oscillations to facilitate speech intelligibility Elife 5 1-17
51 52	602	[31]	Mégevand P, Mercier M R, Groppe D M, Golumbic E Z, Mesgarani N, Beauchamp M S,
53 54	603		Schroeder C E and Mehta A D 2018 Phase resetting in human auditory cortex to visual speech
55 56 57	604	(bioRxiv 1–21
58 59	605	[32]	Ferraro S, Van Ackeren M J, Mai R, Tassi L, Cardinale F, Nigri A, Bruzzone M G, D'Incerti
60	606	7	L, Hartmann T, Weisz N and Collignon O 2020 Stereotactic electroencephalography in

1 2			
3 4	607		humans reveals multisensory signal in early visual and auditory cortices Cortex
5 6	608	[33]	Arnal L H and Giraud A L 2012 Cortical oscillations and sensory predictions Trends Cogn.
7 8	609		Sci. 16 390–8
9 10	610	[34]	Gross J, Hoogenboom N, Thut G, Schyns P, Panzeri S, Belin P and Garrod S 2013 Speech
11 12	611		Rhythms and Multiplexed Oscillatory Sensory Coding in the Human Brain PLoS Biol. 11
13 14	612	[35]	Giraud A L, Kleinschmidt A, Poeppel D, Lund T E, Frackowiak R S J and Laufs H 2007
15 16	613		Endogenous Cortical Rhythms Determine Cerebral Specialization for Speech Perception and
17 18	614		Production Neuron 56 1127–34
19 20 21	615	[36]	Etard O and Reichenbach T 2019 Neural speech tracking in the theta and in the delta
21 22 23 24 25 26 27 28 29 30 31 32	616		frequency band differentially encode clarity and comprehension of speech in noise J.
	617		Neurosci. 39 5750–9
	618	[37]	Keshavarzi M, Kegler M, Kadir S and Reichenbach T 2020 Transcranial alternating current
	619		stimulation in the theta band but not in the delta band modulates the comprehension of
	620		naturalistic speech in noise Neuroimage 210 116557
32 33	621	[38]	Ding N and Simon J Z 2014 Cortical entrainment to continuous speech: Functional roles and
34 35 36 37 38 39 40 41 42 43 44	622		interpretations Front. Hum. Neurosci. 8 1–7
	623	[39]	van Wassenhove V and Grzeczkowski L 2015 Visual-induced expectations modulate auditory
	624		cortical responses Front. Neurosci. 9 1-10
	625	[40]	Simon D M and Wallace M T 2017 Rhythmic Modulation of Entrained Auditory Oscillations
	626		by Visual Inputs Brain Topogr. 30 565–78
45 46	627	[41]	Keil J and Senkowski D 2018 Neural Oscillations Orchestrate Multisensory Processing
47 48	628		Neuroscientist 24 609–26
49 50	629	[42]	Luo H, Liu Z and Poeppel D 2010 Auditory cortex tracks both auditory and visual stimulus
51 52	630		dynamics using low-frequency neuronal phase modulation PLoS Biol. 8 25-6
53 54	631	[43]	Stanley D A, Falchier A Y, Pittman-Polletta B R, Lakatos P, Whittington M A, Schroeder C E
55 56 57	632	(and Kopell N J 2019 Flexible reset and entrainment of delta oscillations in primate primary
58 59	633		auditory cortex: modeling and experiment bioRxiv 1-42
60	634	[44]	Pittman-Polletta B R, Wang Y, Stanley D A, Schroeder C E, Whittington M A and Kopell N J
	· · · · · · · · · · · · · · · · · · ·		

2			
3 4	635		2020 1 Differential contributions of synaptic 2 and intrinsic inhibitory currents to 3 speech
5 6	636		segmentation via flexible 4 phase-locking in neural oscillators bioRxiv 1-22
7 8	637	[45]	Kulkarni A, Ranft J and Hakim V 2020 Synchronization, Stochasticity, and Phase Waves in
9 10 11 12 13 14 15 16	638		Neuronal Networks With Spatially-Structured Connectivity Front. Comput. Neurosci. 14
11 12	639	[46]	Hyafil A, Fontolan L, Kabdebon C, Gutkin B and Giraud A L 2015 Speech encoding by
13 14	640		coupled cortical theta and gamma oscillations <i>Elife</i>
15 16 17	641	[47]	Jadi M P and Sejnowski T J 2014 Regulating cortical oscillations in an inhibition-stabilized
17 18 19	642		network <i>Proc. IEEE</i> 102 830–42
20 21	643	[48]	Mazzoni A, Panzeri S, Logothetis N K and Brunel N 2008 Encoding of naturalistic stimuli by
22 23	644		local field potential spectra in networks of excitatory and inhibitory neurons PLoS Comput.
24 25	645		Biol. 4
26 27	646	[49]	Kegler M and Reichenbach T 2021 Modelling the effects of transcranial alternating current
28 29	647		stimulation on the neural encoding of speech in noise Neuroimage 224 117427
30 31	648	[50]	Virtanen P, Gommers R, Oliphant T E, Haberland M, Reddy T, Cournapeau D, Burovski E,
32 33	649		Peterson P, Weckesser W, Bright J, van der Walt S J, Brett M, Wilson J, Millman K J,
34 35 36	650		Mayorov N, Nelson A R J, Jones E, Kern R, Larson E, Carey C J, Polat İ, Feng Y, Moore E
37 38	651		W, VanderPlas J, Laxalde D, Perktold J, Cimrman R, Henriksen I, Quintero E A, Harris C R,
39 40	652		Archibald A M, Ribeiro A H, Pedregosa F, van Mulbregt P, Vijaykumar A, Bardelli A Pietro,
41 42	653		Rothberg A, Hilboll A, Kloeckner A, Scopatz A, Lee A, Rokem A, Woods C N, Fulton C,
43 44	654		Masson C, Häggström C, Fitzgerald C, Nicholson D A, Hagen D R, Pasechnik D V., Olivetti
45 46	655		E, Martin E, Wieser E, Silva F, Lenders F, Wilhelm F, Young G, Price G A, Ingold G L, Allen
47 48	656		G E, Lee G R, Audren H, Probst I, Dietrich J P, Silterra J, Webber J T, Slavič J, Nothman J,
49 50	657		Buchner J, Kulick J, Schönberger J L, de Miranda Cardoso J V, Reimer J, Harrington J,
51 52	658		Rodríguez J L C, Nunez-Iglesias J, Kuczynski J, Tritz K, Thoma M, Newville M, Kümmerer
53 54 55	659		M, Bolingbroke M, Tartre M, Pak M, Smith N J, Nowaczyk N, Shebanov N, Pavlyk O,
55 56 57	660	(Brodtkorb P A, Lee P, McGibbon R T, Feldbauer R, Lewis S, Tygier S, Sievert S, Vigna S,
58 59	661		Peterson S, More S, et al 2020 SciPy 1.0: fundamental algorithms for scientific computing in
60	662	7	Python Nat. Methods 17 261–72

1 2			
3 4 5	663	[51]	Goodman D F M and Brette R 2009 The brian simulator Front. Neurosci. 3 192–7
5 6	664	[52]	John S. Garofolo, Lamel L F, Fisher W M, Fiscus J G, Pallett D S and Dahlgren N L 1990
7 8	665		Acoustic-Phonetic Continuous Speech Corpus
9 10	666	[53]	Cooke M, Barker J, Cunningham S and Shao X 2006 An audio-visual corpus for speech
11 12	667		perception and automatic speech recognition J. Acoust. Soc. Am. 120 2421–4
13 14	668	[54]	Chi T, Ru P and Shamma S A 2005 Multiresolution spectrotemporal analysis of complex
15 16	669		sounds J. Acoust. Soc. Am. 118 887–906
17 18 10	670	[55]	Pillow J W, Shlens J, Paninski L, Sher A, Litke A M, Chichilnisky E J and Simoncelli E P
20 21	671		2008 Spatio-temporal correlations and visual signalling in a complete neuronal population
22 23	672		Nature 454 995–9
7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 34 35 36 37 38 9 40 41 42 43 44 45 46 47 48 49 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 45 45 46 47 48 46 47 48 47 48 48 48 48 48 48 48 48 48 48	673	[56]	Shi J V., Xu Y and Baraniuk R G 2014 Sparse Bilinear Logistic Regression 1–27
	674	[57]	Victor J D 2005 Spike train metrics Curr. Opin. Neurobiol. 15 585–92
	675	[58]	Wilcoxon F 1946 Individual comparisons of grouped data by ranking methods. J. Econ.
	676		Entomol. 39 269
	677	[59]	Benjamini Y and Hochberg Y 1995 Controlling the False Discovery Rate: A Practical and
	678		Powerful Approach to Multiple Testing J. R. Stat. Soc. Ser. B 57 289-300
	679	[60]	Krischer L, Megies T, Barsch R, Beyreuther M, Lecocq T, Caudron C and Wassermann J
	680		2015 ObsPy: A bridge for seismology into the scientific Python ecosystem Comput. Sci.
	681		Discov. 8 0–17
	682	[61]	German-Sallo Z and German-Sallo M 2017 Multiscale Analysing Methods in
	683		Electrocardiogram Signal Processing and Interpretation Procedia Eng. 181 583-7
47 48	684	[62]	Pack C C, Born R T and Livingstone M S 2003 Two-dimensional substructure of stereo and
49 50	685		motion interactions in macaque visual cortex Neuron 37 525-35
51 52	686	[63]	Grant K W and Greenberg S 2001 Speech Intelligibility Derived from Asynchronous
53 54	687		Processing of Auditory-Visual Information Proc. Conf. Audit. Speech Process. 132-7
55 56 57	688	[64]	Brungart D S, Iyer N, Simpson B and Wassenhove V Van 2008 The effects of temporal
57 58 59	689		asynchrony on the intelligibility of accelerated speech Avsp 19-24
60	690	[65]	Ganesan K, Plass J, Beltz A M, Liu Z, Grabowecky M, Suzuki S, Stacey W C, Wasade V S,

2 3 1	691	Towle V L, Tao J X, Wu S, Issa N P and Brang D 2020 Visual speech differentially modulates
5	692	beta, theta, and high gamma bands in auditory cortex <i>bioRxiv</i>
6 7 8 9 10 11 2 13 14 15 16 7 8 9 21 22 32 22 22 22 22 22 22 22 22 22 22 22	693	