# Visual Cortical Tracking of Categorical Speech Features Is Enhanced for Trained Lipreaders

**Zhewei Cao**

Honors Thesis

Department of Brain and Cognitive Sciences,

University of Rochester, Rochester, NY, USA

Spring, 2020

# Abstract

Neuroimaging research has demonstrated that our brain tracks features of observed visual speech during silent lip-reading. Specifically, it has been shown that activity over the occipital scalp can be better decoded when representing speech as a combination of low-level features and categorical speech features, rather than either feature set along, suggesting a tracking of both types of features in visual cortex. However, it remains unclear whether this tracking specifically reflects visible categorical speech features. In addition to visual cortex, it has been shown that auditory cortex is also activated during silent lip-reading, yet not much is known about whether cortical activity is tracking the unheard auditory stimuli in a meaningful way. In the current study, we look into whether silent lip-reading can elicit tracking of categorical visual speech features, as well as the unheard speech envelope in the absence of auditory speech. To do this, we trained participants to be better lip-readers on five audiovisual videos of a speaker, and then tested them on a random selection of the five silent trained videos as well as five silent novel videos of the same speaker while performing a target word detection task. We recorded both behavioral performance and recorded electroencephalography (EEG) data during silent lipreading. Results showed that EEG data from both trained and novel videos saw no clear auditory cortical tracking of the unheard speech signal. However, in visual cortical regions we found enhanced tracking of categorical visual speech features for trained videos over novel videos after regressing out the motion in the videos. Importantly, the extent to which the motion in the videos was tracked by visual regions was the same for both trained and novel videos. With these results, we cast doubt on whether silent lip-reading training does in fact improve cortical representation of the acoustic envelope and found supporting evidence for categorical speech processing at the visual cortex during silent lipreading.

## Introduction

1    Spoken language is a multisensory percept that constitutes the foundation of human interaction. In a face-

2    to-face setting, speech is comprised of not only auditory, but also a visual component. Although the

3    neural processing of auditory speech has been much explored, the visual component of speech however

4    has received much less attention.


5    Studies have shown that the addition of vision enhances speech perception in compromised hearing

6    conditions (Sumby and Pollack, 1954; Ross et al., 2007) and for populations with impaired hearing

7    (Erber, 1971). This benefit is thought to derive from visual speech interacting with auditory speech

8    through two "modes". The first involves vision providing information about the acoustic auditory signal

9    that are hard to detect (Campbell, 2008). For example, it has been shown that from visual speech,

10    specifically mouth width, listeners are able to extract spectro-temporal information of the auditory speech

11    (Plass, 2019). The second mode involves visual speech taking on a correlated role when visual speech

12    information, such as mouth and jaw movement, are temporally correlated with the acoustic auditory

13    information. This redundancy between information from the two modalities has been found to benefit

14    people with normal hearing in optimal hearing conditions (Reisberg *et al.*, 1987; Jiang *et al.*, 2002).


15    A few studies have tried to characterize the neurophysiology underlying visual speech alone, i.e., silent

16    lip-reading. It has previously been shown that cortical tracking to visual-only speech was significantly

17    correlated with subject's lip-reading ability (Crosse *et al.*, 2015). To look into which specific components

18    of visual features the subjects are encoding, there has been evidence for the encoding of higher-level

19    visual speech features in visual cortex (Bernstein & Liebenthal, 2014; O'Sullivan *et al.,* 2017). In one

20    study, it was suggested that combining both low-level visual speech features (e.g., frame-to-frame

21    motion) and higher-level visual features (e.g., visemes – defined as visually similar phonemes) more

22    accurately predicts neural activity over occipital scalp. Yet, further evidence is required to definitively

23    prove that the visual cortex indeed engages in categorical speech processing.

24 Maybe not surprisingly, the auditory cortex might also be playing a role in silent lip-reading -- fMRI

25 literature has suggested the auditory cortex might be activated during silent lipreading (Bernstein *et al.,*

26 2002; Calvert *et al.*, 1997; Pekkola *et al.*, 2005). A model has suggested that visual speech provides

27 temporal cues about the acoustic signal, which could directly project to and affect the sensitivity and

28 activation of auditory cortex (Calvert *et al.*, 1997; Beauchamp *et al.*, 2004). However, little is known

29 about the quantitative nature of cross-modal activation between the visual and auditory cortex during

30 speech perception, specifically, whether the auditory cortex is being activated in a manner that is

31 meaningfully related to the unheard speech.

32 Here, to better understand how visual speech can aid auditory speech perception in various hearing

33 conditions, we trained subjects with natural, continuous audiovisual speech stimuli and recorded how

34 their brain represented the unheard speech when they lipread silent visual speech. Specifically, we

35 hypothesized that when subjects could better lipread, the enhanced lipreading ability would be indexed by

36 (1) enhanced tracking of the unheard acoustic speech envelope in the auditory cortex & (2) improved

37 representation of higher-level visual speech features in the visual cortex. After preliminary analysis, we

38 found that neural signals over visual areas better represented a categorical representation of speech (i.e.,

39 visemes) when participants were previously trained to lip read that speech, suggesting higher-level,

40 speech-specific processing at the visual cortices during silent lipreading, and further illuminating visual

41 cortex's role in audiovisual speech perception.

# Method

## *Subjects*

42    Sixteen native English speakers (11 females; age range: 19–37 years), none of which were trained lip

43    readers, gave written informed consent to participant in the experiment. All participants were right-

44    handed , had self-reported normal hearing, normal or corrected-to-normal vision, and no underlying

45    neurological disease. The study was approved by the Research Subjects Review Board (RSRB) at the

46    University of Rochester.

## *Stimuli*

47    The speech stimuli were drawn from a collection of videos featuring a well-known male speaker,  whose

48    head, shoulders and chest are centered in the frame. The speech was conversational like, and the linguistic

49    content focused on political policy. These speech stimuli were rendered into 60-s 1280 × 720-pixel

50    movies in VideoPad Video Editor (NCH Software). Soundtracks were deleted from the 15 videos which

51    had a frame rate of 30 frames per second.

52    Stimulus presentation and data recording took place in a dark, sound-attenuated room with participants

53    seated at a distance of 70 cm from the visual display. Visual stimuli were presented on a 26″ LCD

54    monitor operating at a refresh rate of 60 Hz.

## *Procedure*

55    Ten different one-minute videos were randomly drawn from the total of 15 videos, five of which selected

56    as training videos.  The experiment was divided into two parts: training and testing.

57    In the first phase, participants were trained to lip-read five randomly selected videos. These videos with

58    intact sound were presented to the participant 10 times each in a random sequence. Subjects were

59    instructed to watch the videos and focus on the lip movements as attentively as they could, while listening

60    to the speech through the provided headphone in the booth. No behavioral testing was carried out at this

61    stage.

62    After the training phase, we recorded EEG while subjects watched silent versions of the five videos

63    learned during the training phase and five novel videos. To refresh participants' memory, each silent

64    trained video was preceded by the same video with sound.  The videos (audiovisual, novel and trained)

65    were randomly placed into a sequence of 15 and were presented to the subjects for a total of 4 iterations

66    during this testing phase.  A demonstration of the testing sequence within an iteration could be found in

67    Figure 1.

**Cond.**  AV, T, N, AV, T, N, N, AV, T .......

**Video**  4, 4, 7, 1, 1, 3, 6, 9, 9 ...........

*Figure 1, Subject's trial sequence in the testing phase. Subjects have four iterations of 15 trials. Within an iteration, there are five audio-visual trained (AV, blue), five silent trained (T, red) and five silent novel (N, green) videos. Each silent trained video (T) always follows its audiovisual version (AV).*

68    To measure how well participants could lip-read videos, we asked participants to rate their subjective

69    intelligibility of the stimuli on a scale from 0-10 and to perform a word detection task. At the beginning

70    of each video's presentation, subjects were given a target word to detect. This word was unique for each

71    trial and we never reused a target word between trial repetitions. Subjects were instructed to press the

72    space bar on the keyboard when they perceived the occurrence of the target word, e.g. reading lip

73    movements in silent trials, or hearing the word in non-silent trials. At the end of each trail, subjects were

74    asked to rate the perceived intelligibility of the watched video on a scale of 0-10, 0 representing that they

75    could understand 0~10% of the video just watched and 10 being that they could understand 90~100%.

76    Subjects were instructed to fixate on the speaker's mouth while minimizing eye blinking and all other

77    motor activity.


*Calculation of behavioral accuracy (D')*

78    The subjects' behavioral responses were recorded by the Presentation software (Version 18.0,

79    Neurobehavioral Systems, Inc., Berkeley, CA), and were used to calculate the d-prime statistic. The d-

80    prime statistic is considered an ideal measure for how readily the signal could be picked up among noise.

81    (Macmillan and Creelman, 2004).  In a signal detection task, it is defined as

$$d' = \frac{\mu_S - \mu_N}{\sqrt{\frac{1}{2}\left(\sigma_S^2 + \sigma_N^2\right)}} \tag{1}$$

82    , with μS and σS signal mean and standard deviation, and μN and σN noise mean and standard deviation.

83    In the case of this experiment, d' serves as a good estimate of the subject's sensitivity to target words in

84    the silent visual speech. Button presses recorded within 2 seconds after target word onset were recorded

85    as hits (signal); and those recorded outside the 2 second timeframe were calculated as false alarms

86    (noise).


*EEG Acquisition and Pre-Processing*

87    During the testing phase, continuous EEG data were acquired from subjects using an ActiveTwo system

88    (BioSemi) from 128 scalp electrodes. The data were low pass filtered online below 134 Hz and digitized

89    at a rate of 512 Hz. We synchronized the EEG to the stimulus via triggers that were sent by an Arduino

90    Uno microcontroller that detected an audio click inserted at the beginning of each soundtrack. EEG signal

91    pre-processing was conducted in MATLAB; the data were bandpass filtered between 0.3 Hz and 15 Hz

92  and re-referenced to the average of all channels. To identify channels with excessive noise, the standard

93  deviation of each channel was compared with that of the surrounding channels in MATLAB. Channels

94  which are two standard deviations away from neighboring channels were deemed contaminated by noise

95  and replaced by spline-interpolating the remaining clean channels with weightings based on their relative

96  scalp location in EEGLAB (Delorme and Makeig, 2004). Each trial's data was down sampled to 64Hz

97  before further analysis.

### *EEG Analysis - Stimulus Feature Extraction*

98  Our overall EEG analysis strategy was based on relating the ongoing dynamics of the recorded EEG to

99  different dynamic features of the audio and visual speech. More particularly, we aimed to fit models of

100  the EEG responses based on different features of the speech, and then to test if those model can

101  successfully predict EEG responses to new stimuli. If some models do a better job than others of

102  predicting EEG, one can then say something about what features of the speech are being reflected in the

103  data and, thus, perhaps, in the brain. A substantial amount of previous research has done this in order to

104  study the hierarchical processing of speech (Di Liberto et al., 2015) and how that processing is affected

105  by attention (Power et al., 2012) and multisensory input (Crosse et al., 2015, 2016).

106  To test our hypotheses about auditory cortical involvement in silent lip-reading and the possibility that

107  visual cortex might be processing categorical features of visual speech, we represented the speech using

108  several different representations.

109  ***Envelopes*** Here we use the speech envelope as an estimate of the unheard acoustic speech, as

110  demonstrated in previous studies (O'Sullivan *et al.,* 2017). The broadband amplitude envelope

111  representation was obtained by bandpass filtering the speech signal into 256 logarithmically spaced

112  frequency bands between 80 Hz and 3000 Hz using a gammachirp filter bank (Irino and Patterson, 2006).

113    The envelope at each of the 256 frequency bands was calculated using a Hilbert transform, and the

114    broadband envelope was obtained by averaging over the 256 narrowband envelopes.

115    ***Frame-to-frame Motion*** To accurately capture the motion of pixels in the videos, a frame-to-

116    frame motion was calculated. For each frame, a matrix of motion vectors was calculated using an

117    ''Adaptive Rood Pattern Search'' block matching algorithm (Barjatya, 2004). Through pooling all motion

118    vector lengths in the frame (Bartels et al., 2008), a global motion vector was obtained. This vector was

119    resampled to 64 Hz to match the sampling rate of the EEG data.

120    ***Phoneme/Viseme***  Groupings of phonemes have been identified to be ambiguous and

121    confusable when presented visually during identification tasks (Woodward and Barber, 1960; Fisher,

122    1968). Therefore, each of these groups of visually confusable phonemes can be seen as the building

123    blocks of visual speech, i.e. visemes. To derive a viseme representation from our videos, a phonemic

124    representation from Di Liberto et al. (2015) was used. Using methods defined in Auer and Bernstein

125    (1997), the phonemic representation was converted to visemes. The phoneme- to-viseme transformation

126    means that timing of our viseme representation is actually synchronized with the acoustic boundaries

127    rather than the visual.

### EEG Analysis - Temporal Response Function Estimation

128    In order to relate the continuous EEG to the various visual speech representations introduced above, we

129    used the temporal response function (TRF), a regression analysis that enables mapping between EEG and

130    features. A TRF can be thought of as a filter that describes how a particular stimulus feature (e.g., the

131    acoustic envelope) is transformed into the continuous EEG at each channel and is analogous to the event

132    related potential. So if $s(t)$ represents the stimulus feature at time $t$, the EEG response at channel $n$, $r(t, n)$,

133    can be modeled as a convolution with a to-be-estimated TRF, $w(\tau, n)$.

$$r(t, n) = \sum_{\tau = T_{\min}}^{T_{\max}} w(\tau, n)s(t - \tau) + \varepsilon(t, n), \qquad (2)$$

134    where ε(t, n) stands for responses unexplained by the model at each channel. Considering that it takes

135    several tens of milliseconds for the effect of stimulus to be detectable in the responses, the TRF is

136    calculated between stimulus and response across different lengths of time-lags. In the current study, we fit

137    TRFs for each 60-s trial using ridge regression expressed in the following matrix form:

$$w = (S^T S + \lambda I)^{-1} S^T r, \qquad (3)$$

138    where λ is the ridge parameter, chosen to optimize the stimulus- response mapping which also provides

139    regularization and prevents overfitting, S is a matrix containing a time series of stimulus samples for the

140    window of interest (i.e., the lagged time series), $r$ is a matrix of the 128-channel neural response data, and

141    $I$ is the identity matrix.  The TRF was computed using a custom-built toolbox in MATLAB (Crosse *et al.*,

142    2016).

### *EEG Prediction and Model Evaluation*

143    With this TRF modeling approach, we set out to assess how some of the abovementioned speech features

144    was being encoded in neural signals using forward prediction, and how some of the univariate speech

145    features were represented in the EEG signals with backward reconstruction. We fit TRFs describing the

146    mapping between the speech features and the EEG. Then, using leave-one-out cross-validation, we

147    examine how well the left-out EEG data could be predicted using the different models. If the EEG data

148    could be predicted with accuracy using a particular model, we can suggest that the EEG is reflecting the

149    encoding of that particular feature or set of features. Because we had 20 trials for each subject, leave-one-

150    out cross-validation meant that each TRF was fit to the data from 19 trials and then the average TRF

151    across these 19 trials was used to predict the EEG in the remaining trial (Crosse *et al.,* 2016; O'Sullivan *et*

152    *al.,* 2017).

153   Prediction accuracy was measured by calculating Pearson's (r) linear correlation coefficient between the

154   predicted and original EEG responses at each electrode channel.  This procedure is explained in more

155   detail in Crosse et al. (2016).  In the case of this study, when two features A and B are related, we

156   attempted to partial out the influence of one from another – For every trial, after we generated predicted

157   EEG from one speech feature A, the predicted signal is subtracted from the true EEG responses. This

158   gives us EEG that has the contribution from feature A partialled out. We then try to predict this residual

159   EEG using feature B. This tells us how well feature B is represented in the EEG after accounting for the

160   contribution of feature A.

# *Results*

## *Training improves lip-reading ability*

161   All sixteen participants performed a target word detection task while EEG data was acquired from them.

162   To investigate the effect of training on the task performance, we compared the self-reported measure of

163   intelligibility, as well as sensitivity to the target word during the silent lipreading trials.

164   The training's effect can be measured by both the subject's self-reported measures of intelligibility after

165   each trial, as well as their behavioral performances during the trial. Subjects rated videos that they were

166   trained on as significantly more intelligible (Figure 2a, p = 0.0059, Wilcoxon test). Consistent with their

167   self-report, subjects' performance in target word detection task also improved on trained videos (Figure

168   2b). The condition had a significant effect on d-prime (p = 0.0076, Wilcoxon test).



*Figure 2, Subject's behavioral responses during the lip-reading task: A) Subject's rating of the silent speech videos for conditions of novel and trained ; B) Subject's performance of target word detection for the silent speech videos for conditions of novel and trained*

### *Unheard Speech Envelope reconstruction is not improved after training*

169    To look into whether the training had an impact on the cortical representation of unheard speech, we

170    reconstructed an estimate of the unheard speech envelope with the collected EEG signal for trained and

171    novel trials (Figure 3a). We found that the unheard envelope could not be more accurately reconstructed

172    within each subject in the trained condition compared with the novel condition (p=0.40, Wilcoxon test).

173    Following this finding, we set out to examine whether regions of the scalp were contributing differently to

174    the two conditions -- we predicted data collected from each EEG channel with the acoustic envelope and

175    visualized the differences between them (Figure 3b). Although there seems to be differences between the

176    two conditions in the frontal central scalp, signals collected from the frontal-central channels found no

177    difference between novel and trained conditions (Figure 3c).



*Figure 3a, the reconstruction accuracies of unheard speech envelope for novel and trained conditions (n=16, p=0.40, Wilcoxon)*

*Figure 3b Topographic maps of prediction accuracies using the unheard speech envelope*
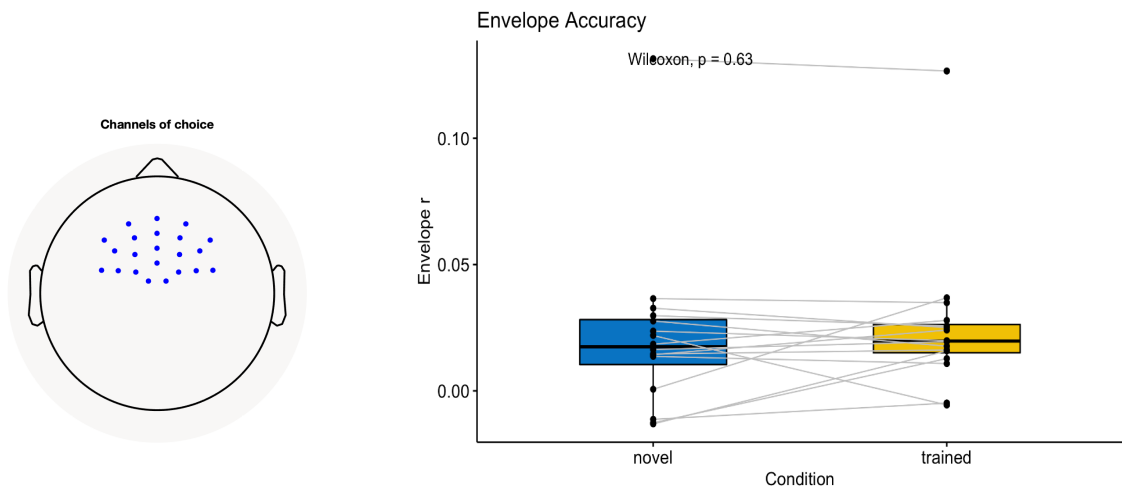


*Figure 3c Left: 22 channels chosen over the frontal-central region*

*Right: Mean prediction accuracy accuracies using the acoustic speech envelope over the chosen channels for each subject in both novel and trained conditions*

## *Encoding of Frame-to-Frame Motion*

178    The frame to frame motion stimulus used in the current study includes both local and global motion in the

179    visual scene (Bartels et al., 2008, O'Sullivan et al., 2017). While some of the global and local motion

180    components may be tangential to speech (e.g., movement of speaker relative to the background), some of

181    them are related (e.g. jaw movements). In this way, the frame-to-frame vector provides an estimate of

182    general low-level visual information received by the visual cortices.

183    To look into whether training has an impact on the cortical responses to the silent speech videos, we first

184    examined how the frame-to-frame motion is encoded in both trained and novel conditions (Figure 4a). We

14

185    showed an enhancement in motion tracking for trained vs novel conditions particularly in the visual

186    regions. However, to account for the fact that motion and visemes are correlated with each other (Files *et*

187    *al*., 2015), and to examine frame to frame motion's unique contribution to visual speech perception, we

188    first regressed out visemes, and then related the residual EEG to the motion vectors (Figure 4b). We then

189    chose 25 channels at the occipital area and found no difference between trained and novel conditions

190    (Figure 4c, p=0.38, Wilcoxon test).



*Figure 4a Topographic maps of prediction accuracies using frame to frame motion in the unheard speech*



*Figure 4b Topographic maps of prediction accuracies using frame to frame motion in the unheard speech after regressing out*

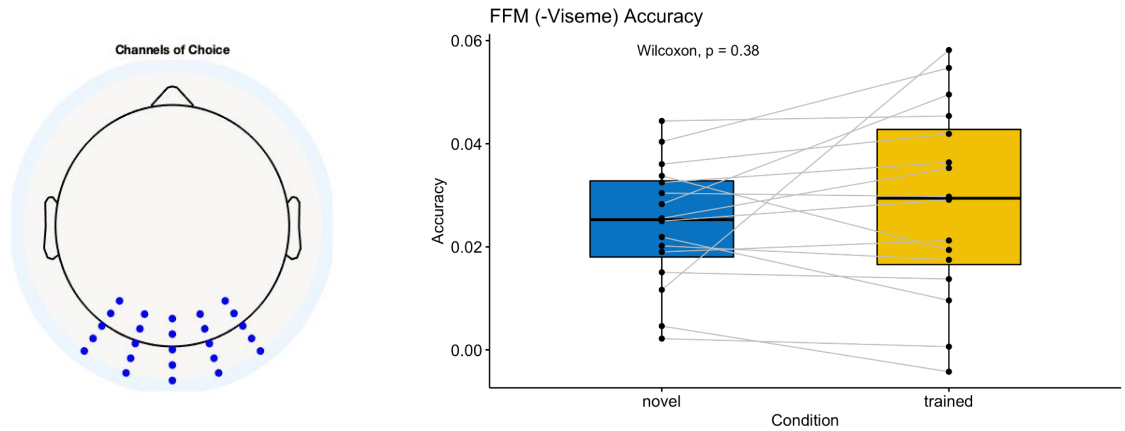*potential contributions from visual phonetic features*

15

*Figure 4c Left:25 channels chosen over the occipital region*

*Right: Mean prediction accuracy accuracies using frame-to-frame motion while excluding visual phonetic features over the chosen channels for each subject in both novel and trained conditions*

## *Encoding of Visemes*

191    It has been suggested that EEG may be reflecting categorical speech processing beyond encoding the

192    onset and offset of speech, yet more evidence is needed for whether speech-specific, categorical-viseme

193    level processing were ongoing at the visual cortices. We reasoned that, if such processing can happen in

194    visual areas, it should be specifically enhanced during successful lipreading. To test this, we wanted to

195    examine how training participants to lipread certain videos would affect the encoding of visemes. In both

196    auditory (Di Liberto et al., 2015) and visual speech (O'Sullivan et al., 2017; Hauswald et al., 2018), the

197    relationship between EEG and phoneme representations of speech has been explored.

198    We found that the higher prediction accuracies for visemes in the visual regions for trained conditions

199    (Figure 5a). Again, to assess the unique contribution of viseme to visual speech perception, we regressed

200    out frame-to-frame motion and related the residual EEG data to the viseme vectors (Figure 5b). To

201    quantify the enhanced tracking of visemes only at the occipital area, we examined 25 channels over the

202      occipital region and found significant improvement in viseme representation for trained videos compared
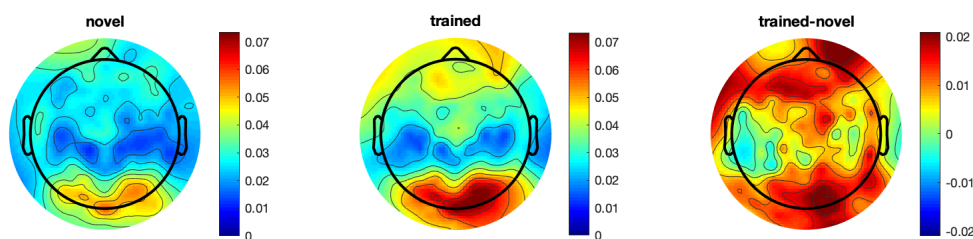
203      with novel (Figure 5c, p=0.0092, Wilcoxon test).



*Figure 5a Topographic maps of prediction accuracies using visual phonetic features (i.e. visemes) in the unheard speech*
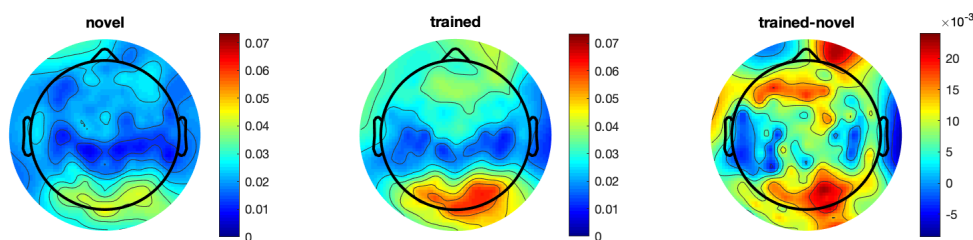


*Figure 5b Topographic maps of prediction accuracies using visual phonetic features (i.e. visemes) in the unheard speech, after regressing out the influence from low level visual features (i.e. frame-to-frame motion).*
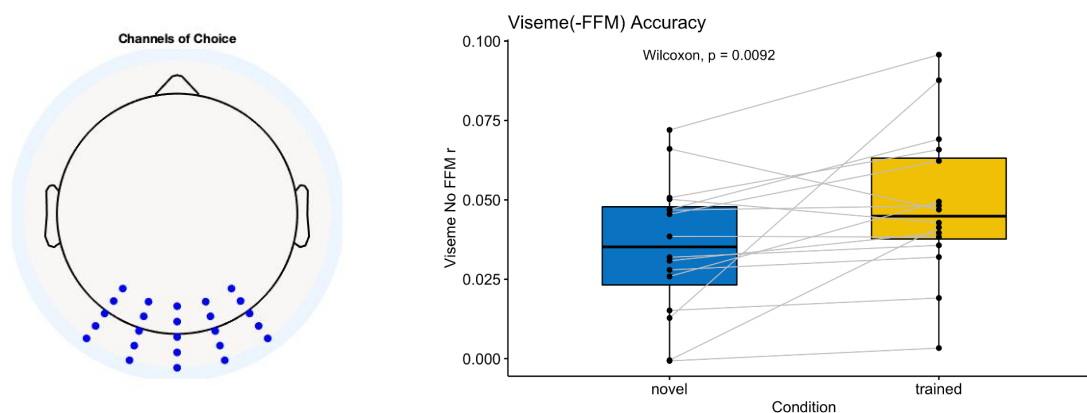


*Figure 5c  Left:25 channels chosen over the occipital region*

*Right: Mean prediction accuracy accuracies using visual phonetic features (i.e. visemes) while excluding frame-to-frame motion over the chosen channels for each subject in both novel and trained conditions*

## *Discussion*

204   In the current study, we trained subjects to lipread with audiovisual videos. Our results show that this

205   training was successful. We demonstrate that during lipreading of trained videos, neural signals are

206   representing categorical representations of speech more robustly. This suggests that the visual cortex is

207   engaged in processing of categorical visual speech features, besides encoding of physical stimulus

208   dynamics.

### *Training didn't improve tracking of the unheard speech envelope*

209   Our finding that cortical representation of the unheard speech envelope wasn't enhanced when subjects

210   were able to better lipread (Figure 3) seems to be at odds with previous findings that showed cortical

211   tracking of unheard acoustic information (Bourguignon et al. 2020). Bourguignon et al. showed that

212   auditory cortices are entraining to frequencies below 1 Hz, which matches with rhythmicity of phrases

213   and sentences. However, they failed to observe a similar tracking for frequency between 4-8 Hz, the

214   syllable rate in the speech. Contrary to previous findings of enhanced tracking for frequency between 4-8

215   Hz in the auditory cortex for intelligible auditory speech than unintelligible (Ahissar et al., 2001; Luo and

216   Poeppel, 2007; Peelle et al., 2013),  the current study like Bourguignon and colleagues, also failed to find

217   similar differences in tracking of silent visual speech between more intelligible (trained) and less

218   intelligible (novel) video conditions.

### *Representation of categorical visual speech features is enhanced with improved lipreading abilities*

219   We found no difference in tracking of low-level visual features in the occipital area between trained and

220   novel conditions, suggesting that this said tracking is not speech specific. This finding is consistent with

221   previous work in the field(Hauswald *et al.* 2018). Hauswald and colleagues found no difference in

222    tracking of low-level visual features when silent visual speech videos were played forward and backward,

223    also implying the non-speech-specific nature of tracking of motion in the visual cortices.

224    We also found supporting evidence for the encoding of categorical speech-specific visual features in the

225    occipital area. This is also in line with previous work, which has also shown that EEG signals in the

226    occipital area can be better decoded with a model with categorical visual speech features than one without

227    (O'Sullivan et al., 2017). Further, other studies have identified posterior superior temporal sulcus (pSTS)

228    and the posterior medial temporal gyrus (MTG) as loci for visuo-phonological processing (Bernstein *et*

229    *al.*,2011), and preliminary results from the current study supports this hypothesis.


### *Limitations and Future Directions*

230    It is important to consider some limitations of the current study. First, the training paradigm of the current

231    study might not be entirely effective due to the length of the training stimuli. For the normal hearing

232    population, 60 seconds of audiovisual stimuli is a lot to memorize and lipreading is a difficult task after

233    all. Second, in the vein of poor training outcomes, for the target word detection task, some subjects might

234    be remembering visual markers of the target word, instead of truly learning to lipread. This would've led

235    to a higher sensitivity to the target word despite showing similar enhancement in cortical tracking.

236    Importantly, however, if this was true, it would only decrease the effect size of any visual speech

237    processing. As such, our increased measure of viseme processing for trained speech is likely a lower

238    bound.

239    The current study could benefit from a number of further statistical analysis on the current data.  First of

240    all, it'd be interesting to explore cortical representations of specific frequency-bands of the unheard

241    auditory stimuli to answer questions like whether the auditory cortex is better at tracking <1 Hz frequency

242    during trained trials. Second of all, it would also be of interest to examine individual differences in the

243  collected data. Such analysis could shed light on the relationships between behavioral measures, e.g.

244  improvement in intelligibility and that in lipreading abilities; as well as correlation between behavioral

245  and neurophysiological measures, e.g. improvement in lipreading abilities and the higher entrainment to

246  stimuli in the cortex.

## *Acknowledgements  & Contributions*

# Reference

Ahissar E, Nagarajan S, Ahissar M, Protopapas A, Mahncke H, Merzenich MM (2001) Speech comprehension is correlated with temporal response patterns recorded from auditory cortex. *Proc Natl Acad Sci U S A* 98:13367–13372.

Barjatya, A. (2004). Block matching algorithms for motion estimation. *IEEE Trans. Evol. Comput*. 8, 225–239.

Bartels, A., Zeki, S., and Logothetis, N. K. (2008). Natural vision reveals regional specialization to local motion and to contrast-Invariant, global flow in the human brain. *Cereb. Cortex* 18, 705–717. doi: 10.1093/cercor/ bhm107

Bernstein, L. E., Auer, E. T. J., Moore, J. K., Ponton, C. W., Don, M., & Singh, M. (2002). Visual speech perception without primary auditory cortex activation. *NeuroReport*, 13(3), 311–315.

Bernstein, L. E., & Liebenthal, E. (2014). Neural pathways for visual speech perception. *Frontiers in Neuroscience*, 8, 386. doi:10.3389/fnins.2014.00386

Beauchamp, M. S., Lee, K. E., Argall, B. D., & Martin, A. (2004). Integration of Auditory and Visual Information about Objects in Superior Temporal Sulcus. *Neuron*, 41(5), 809–823. https://doi.org/10.1016/S0896-6273(04)00070-4

Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2020). Lip-Reading Enables the Brain to Synthesize Auditory Features of Unknown Silent Speech. *Journal of Neuroscience*, *40*(5), 1053–1065.

Calvert, G. A., dullmore, E. T., Brammer, M. J., Campbell, R., Williams, S. C., Mcguire, P. K., et al. (1997). Activation of auditory cortex during silent lipreading. *Science* 276, 593–596. doi: 10.1126/science.276.5312.593

Campbell, R. (2008). The processing of audio-visual speech: empirical and neural bases. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 363, 1001–1010. doi: 10.1098/rstb. 2007.2155

Chandrasekaran C, Trubanova A, Stillittano S, Caplier A, Ghazanfar AA (2009) The natural statistics of audiovisual speech. PLoS Comput Biol 5:e1000436.

Crosse, M. J., Butler, J. S., & Lalor, E. C. (2015). Congruent visual speech enhances cortical entrainment to continuous auditory speech in noise-free conditions. *Journal of Neuroscience*, *35*(42), 14195-14204.

Crosse, M. J., Di Liberto, G. M., Bednar, A., and Lalor, E. C. (2016). The multivariate temporal response function (mTRF) toolbox: a MATLAB toolbox for relating neural signals to continuous stimuli. *Front. Hum. Neurosci.* 10:604. doi: 10.3389/fnhum.2016.00604

Delorme, A., and Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J. Neurosci. Methods* 134, 9–21. doi: 10.1016/j.jneumeth.2003. 10.009

Di Liberto, G. M., O'Sullivan, J. A., and Lalor, E. C. (2015). Low-frequency cortical entrainment to speech reflects phoneme-level processing. *Curr. Biol.* 25, 2457–2465. doi: 10.1016/j.cub.2015.08.030

Erber, N. P. (1971). Auditory and audiovisual reception of words in low-frequency noise by children with normal hearing and by children with impaired hearing. *J. Speech Hear. Res.* 14, 496–512. doi: 10.1044/jshr.1403.496

Files, B. T., Tjan, B. S., Jiang, J., & Bernstein, L. E. (2015). Visual speech discrimination and identification of natural and synthetic consonant stimuli. *Frontiers in Psychology*, 6. https://doi.org/10.3389/fpsyg.2015.00878

Fisher, C. G. (1968). Confusions among visually perceived consonants. *J. Speech Hear. Res*. 11, 796–804. doi: 10.1044/jshr.1104.796

Goncalves, N. R., Whelan, R., Foxe, J. J., and Lalor, E. C. (2014). Towards obtaining spatiotemporally precise responses to continuous sensory stimuli in humans: a general linear modeling approach to EEG. *Neuroimage* 97, 196–205. doi: 10. 1016/j.neuroimage.2014.04.012

Hauswald, A., Lithari, C., Collignon, O., Leonardelli, E., & Weisz, N. (2018). A Visual Cortical Network for Deriving Phonological Information from Intelligible Lip Movements. *Current Biology: CB*, 28(9), 1453-1459.e3. https://doi.org/10.1016/j.cub.2018.03.044

Jiang J, Alwan A, Keating PA, Auer ET, Bernstein LE (2002) On the relationship between face movements, tongue movements, and speech acoustics. EURASIP J Appl Signal Process 11:1174 –1188

Lazard DS, Giraud A-L (2017) Faster phonological processing and right occipito-temporal coupling in deaf adults signal poor cochlear implant outcome. *Nat Commun* 8:14872.

Luo H, Poeppel D (2007) Phase patterns of neuronal responses reliably discriminate speech in human auditory cortex. *Neuron* 54:1001–1010.

Macmillan, N. A., & Creelman, C. D. (2004). Detection theory: A user's guide. Psychology press.

McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, *264*(5588), 746.

O'Sullivan, A. E., Crosse, M. J., Di Liberto, G. M., & Lalor, E. C. (2017). Visual cortical entrainment to motion and categorical speech features during silent lipreading. *Frontiers in Human Neuroscience*, 10, 679. doi:10.3389/fnhum.2016.00679

Peelle JE, Gross J, Davis MH (2013) Phase-locked responses to speech in human auditory cortex are enhanced during comprehension. *Cereb Cortex* 23:1378–1387.

Peelle, J. E., and Sommers, M. S. (2015). Prediction and constraint in audiovisual speech perception. *Cortex* 68, 169–181. doi: 10.1016/j.cortex.2015.03.006

Pekkola, J., Ojanen, V., Autti, T., Jääskeläinen, I. P., Möttönen, R., Tarkiainen, A., et al. (2005). Primary auditory cortex activation by visual speech: an fMRI study at 3 T. *Neuroreport* 16, 125–128. doi: 10.1097/00001756-200502080-00010

Plass, J., Brang, D., Suzuki, S., & Grabowecky, M. (2019, May 20). Vision Perceptually Restores Auditory Spectral Dynamics in Speech.

Power AJ, Foxe JJ, Forde EJ, Reilly RB, Lalor EC (2012) At what time is the cocktail party? A late locus of selective attention to natural speech. *Eur J Neurosci* 35:1497–1503.

Reisberg D, McLean J, Goldfield A (1987) Easy to hear but hard to under- stand: a lip-reading advantage with intact auditory stimuli. In: Hearing by eye: the psychology of lip-reading (Dodd B, Campbell R, eds), pp 97–114. Hillsdale, NJ: Erlbaum.

Ross, L. A., Saint-Amour, D., Leavitt, V. M., Javitt, D. C., & Foxe, J. J. (2006). Do you see what I am saying? Exploring visual enhancement of speech comprehension in noisy environments. *Cerebral cortex*, *17*(5), 1147-1153.

Simpson, A. J.; Fitter, M. J. (1973). "What is the best index of detectability?". Psychological Bulletin. 80 (6): 481–488.

Sumby, W. H., and Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *J. Acous. Soc. Am*. 26, 212–215. doi: 10.1121/1.1907309

Woodward, M. F., and Barber, C. G. (1960). Phoneme perception in lipreading. *J. Speech Hear. Res*. 3, 212–222. doi: 10.1044/jshr.0303.212