# Independent mechanisms of temporal and linguistic cue correspondence benefiting audiovisual speech processing

Sara Fiscella,[1,2] Madeline S Cappelloni,[2,3] Ross K Maddox,[2,3,4,5]

[1] *Brain and Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA*

[2] *Del Monte Institute for Neuroscience, University of Rochester, Rochester, NY 14627, USA*

[3] *Biomedical Engineering, University of Rochester, Rochester, NY 14627, USA*

[4] *Neuroscience, University of Rochester, Rochester, NY 14627, USA*

[5] *Center for Visual Science, University of Rochester, Rochester, NY 14627, USA*

**ABSTRACT**

When listening is difficult, seeing the face of the talker aids speech comprehension. Faces carry both temporal (low-level physical correspondence of mouth movement and auditory speech) and linguistic (learned physical correspondences of mouth shape (viseme) and speech sound (phoneme)) cues. Listeners participated in two experiments investigating how these cues may be used to process sentences when maskers are present. In Experiment I, faces were rotated to disrupt linguistic but not temporal cue correspondence. Listeners suffered a deficit in speech comprehension when the faces were rotated, indicating that visemes are processed in a rotation-dependent manner, and that linguistic cues aid comprehension. In Experiment II, listeners were asked to detect pitch modulation in the target speech with upright and inverted faces that either matched the target or masker speech such that performance differences could be explained by binding, an early multisensory integration mechanism distinct from traditional late integration. Performance in this task replicated previous findings that temporal integration induces binding, but there was no behavioral evidence for a role of linguistic cues in binding. Together these experiments point to temporal cues providing a speech processing benefit through binding and linguistic cues providing a benefit through late integration.

## I.      INTRODUCTION

While having a conversation may be easy in quiet environments, noisy environments can render listening a challenging task. Though many of us can listen to our conversation partner with minimal conscious effort, the neural processes by which we accomplish this auditory feat are complex, rely on many stimulus cues, and are poorly understood.

When auditory information is insufficient or difficult to process, visual cues help us listen. Seeing the face of a talker greatly improves our ability to comprehend their speech (Arnold and Hill, 2001; Reisberg et al., 1987). Speaking faces carry both temporal and linguistic cues that are congruent with auditory speech. Specifically, the movements of the talker's mouth and surrounding areas are temporally coherent with the unique amplitude envelope of the auditory stream of interest, and the shapes the mouth makes are linguistically congruent with the speech sounds produced. Temporal information is an inherently physical cue and constrained by the dynamics of speech production. The time correlation of mouth movements and speech can help listeners pair the relevant auditory and visual streams (Maddox et al., 2015) or even reduce masking effects of competing auditory streams (Grant and Bernstein, 2019). Linguistic information is provided by the link between specific mouth shapes, called visemes, and the phonemes that they generate. Unlike pure temporal coherence, the link between visemes and phonemes is a learned prior that relies on a listener's experience with language. The underlying mechanisms that may allow multisensory linguistic information to help us listen are not known.

Many researchers have turned to the McGurk effect to demonstrate the effect of visual linguistic cues on auditory perception. The McGurk effect occurs when observers are concurrently presented with an auditory syllable (e.g. "ba") and a face which either matches the auditory syllable (e.g. "ba") or a different syllable (e.g. "ga"). Even though subjects are presented with the same auditory syllable in both visual conditions they often report hearing a fused syllable (e.g. "da") when the face and

2

25  auditory speech do not match (Mcgurk and Macdonald, 1976). In order to better understand how

26  the brain processes the linguistic cues associated with the face, studies have looked at the effects of

27  inverting the face. These studies have found that when the face is inverted listeners less accurately

28  identify syllables in the visual alone ("lipreading") condition, and a higher proportion of people

29  report hearing the auditory syllable than the fused syllable in a multisensory (McGurk) condition

30  (Massaro and Cohen, 1996; Ujiie et al., 2018). Despite the findings that inversion of the face can

31  disrupt the processing that underlies the McGurk effect, the monosyllabic stimuli involved do not

32  well model the demands of listening to speech in noise. It is still unclear whether these linguistic

33  cues carry the same perceptual weight in continuous speech.

34      Though temporal and linguistic cues both can contribute to speech comprehension, they may

35  contribute in very different ways or at different stages of the multisensory perception process.

36  Multisensory integration has been traditionally thought of as a Bayesian combination of unisensory

37  information just prior to perceptual decision making ("late integration") (Körding et al., 2007), but

38  recent work has pointed towards an earlier stage of multisensory integration known as "binding" or

39  "early integration" (Atilgan et al., 2018; Bizley et al., 2016; Lee et al., 2019). Binding occurs when an

40  auditory and visual stream are combined by the brain into a single perceptual object. Binding affects

41  the encoding of an audiovisual object, whereas many of the effects of multisensory integration can

42  be explained by a later decision bias (Bizley et al., 2016). When the brain forms a perceptual object, it

43  can then allocate object-based attention (Shinn-Cunningham, 2008). By attending an object, all

44  features are automatically enhanced, even orthogonal features that are not comodulated (Lee et al.,

45  2019). This can both occur within a modality (i.e., multiple visual features combine to form a single

46  visual object (Blaser et al., 2000)) or across modalities (i.e., a visual feature is bound with an auditory

47  feature to form a multisensory object (Maddox et al., 2015)).

3

48      Despite potentially having independent neural underpinnings, studies often fail to distinguish

49      between binding and late integration. Binding can be tested in a task that requires the listener to

50      attend two streams to complete a dual task. If the stimuli in the two streams are bound, they can

51      attend to the combined object and improve their performance instead of having to divide their

52      attention to complete the task (Bizley et al., 2016). For binding to occur, there must be some

53      compelling relationship between the object's features. Possible relationships between features, which

54      may or may not contribute to binding, may roughly be divided into several categories: low-level

55      physical correspondence, semantic congruence, and learned physical correspondence. Low-level

56      physical correspondence involves fundamental relationships of stimuli such as temporal coherence,

57      which is known to induce binding (Maddox et al., 2015), and spatial congruence. Semantic

58      congruence encompasses higher-level learned relationships between stimuli that are commonly

59      paired in the natural world, such as the image and sound of a dog barking, and are unlikely to induce

60      binding due to the significant high-level processing required to make these associations. Learned

61      physical relationships are similar to semantic congruence in the sense that they must be learned

62      through observation of the natural world and similar to low-level physical correspondence in that

63      the physical relationship of the stimuli is inherent to their production, such as visemes and

64      phonemes in which mouth shapes are innately connected to the sounds they produce but are only

65      known to be related by someone with experience of talkers. Here we investigate whether the learned

66      physical relationships of natural speech induce binding.

67      In Experiment I, we sought to determine if rotating the face of a talker disrupts linguistic cues in

68      a behaviorally relevant task and therefore provide some insight into how these cues contribute to

69      listening in noisy environments. We engaged listeners in a speech in noise task with rotated videos

70      of the target talker to determine how their performance was affected by disrupted cues. We found

71    that rotating the face hindered their speech comprehension, suggesting that the information carried

72    by the face was indeed disrupted by the rotation.

73        We tested whether linguistic cues could induce binding in Experiment II. Given that the face

74    inversion disrupted linguistic cues, we looked for differences in a multisensory selective attention

75    task that might suggest differences in binding. Here we asked listeners to detect auditory pitch

76    modulations and visual events in target stimuli while ignoring maskers. Binding is a likely mechanism

77    to explain any improvement in performance when the listener saw the target's face relative to the

78    masker's face. We found that although there was a clear advantage to seeing the target's face in all

79    conditions, there was no effect of face rotation, suggesting that the disrupted linguistic cues did not

80    impair binding.

81        We ultimately find that while linguistic cues are important for listening to speech in noisy

82    environments, this is likely due to late integration rather than binding.
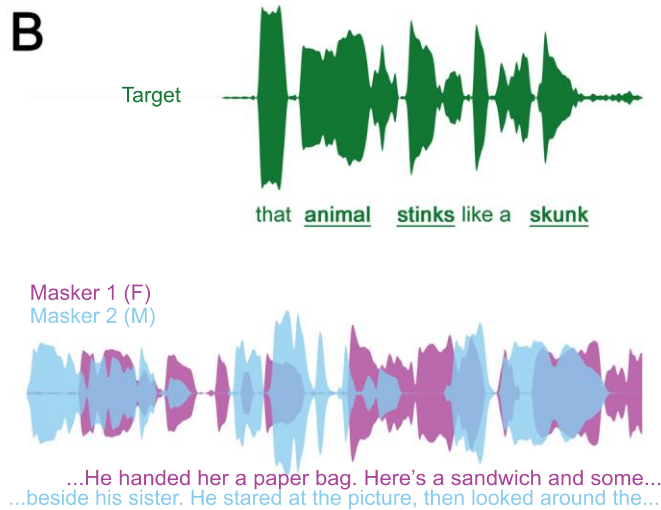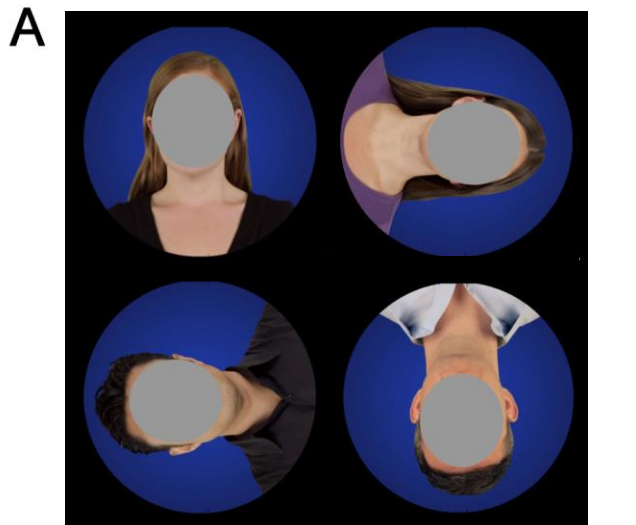
## II.    EXPERIMENT I

### A.  Methods

83

84

85    Participants performed a speech in noise comprehension task.

86    a.  *Participants.* Fifteen participants (12 female, 3 male) age 18−26 (mean 22.6) had normal hearing

87    (20 dB HL or better at octave frequencies from 500 Hz to 8000 Hz), self-reported normal or

88    corrected-to-normal vision, and spoke English as their primary language. Participants gave

89    written consent and were paid for their participation. All protocols were approved by the

90    University of Rochester Research Subjects Review Board.

91    b.  *Stimuli.* The visual and auditory stimuli were selected from the STeVi corpus. The corpus

92    includes video recordings of four native English speakers (two male and two female) saying 200

93    high probability sentences that contained three to five keywords each (e.g. "The scarf was made

94    of shiny silk.") (STeVi Speech Test Video Corpus, n.d.). We rotated the videos by 0°, ±90°, or

95    180° in each trial. In some trials, still frames from the beginning of the unrotated videos were

96    used instead of the dynamic faces. To ensure that all videos were the same size and aspect ratio,

97    we drew a circle for each trial around the face of the speaker and its dark blue background,

98    leaving the background outside of the circle black. In addition, the mouths of the talkers were

99    centered on the screen (Figure 1). We used high probability sentences due to previous research

100   showing a benefit of semantic context in understanding speech in noisy environments (Van

101   Engen et al., 2014). Due to an error in the experiment code, the videos were played at 29.04

102   frames per second (fps) instead of their native original 29.97 fps, which led to a small offset that

103   increased throughout the trial with the biggest offset would be at the end of the longest

104   sentence—a delay of 105 ms in the worst case. On average the delay was 37 ms, about one

105   frame of video, and well within the audio-visual temporal binding window (Stevenson et al.,

106   2012).

6

107    There were also two auditory masker streams comprised of natural speech from American

108    English audiobooks, The Alchemyst (Scott, 2008) (male narrator) and A Wrinkle in Time

109    (L'Engle, 2006) (female narrator). Audio was edited to remove silent pauses longer than 0.5

110    seconds. The masker stimuli were each presented at 60 dB SPL. Then target stimuli were

111    presented with a signal-to-noise ratio (SNR) of 0, −3, or −6 decibels (dB).



112

113    Figure 1: (Color Online). A. Images of the four talkers (faces covered to protect identity of the

114        actors), with 0° (top left), +90° (top right), −90° (lower left), and 180° (lower right) rotation.

115        The same circular mask is applied to all rotations with the mouth centered. B. Envelopes and

116        transcription of auditory stimuli in a given trial. The target sentence began after a 2 second delay

117        (top, green) while two maskers (bottom, female narrator: purple, male narrator: pale blue) played

118        continuously throughout the trial. Subjects had to type each keyword (bold and underlined) in

119        the target sentence while ignoring the maskers in order to receive credit.

120    c.  *Procedure.* Subjects were seated in a dark soundproof booth in front of a 24 inch BenQ monitor,

121        with their nose lined up approximately with the center of the screen and a 50 centimeter viewing

122        distance. Sounds were presented via ER-2 insert earphones (Etymotic Research, Elk Grove

123        Village, Il). Subjects were given a standard keyboard to type in their responses.

124        Subjects were required to pass a training module before beginning the experiment. Training

125        began with two trials without maskers, followed by three trials of with maskers and an SNR of 0

126        dB, and lastly three trials with SNR of −3 dB. Participants responded by typing the sentence

127        they heard after each trial with no capitalization or punctuation. For training, where accuracy

128        needed to be judged in real time, responses in a given trial were scored as correct if the sequence

129        of letters of the entered keywords were at least 80% correct. Subjects were given two chances to

130        pass the training, which required correct responses in both of the trials without background

131        noise and two out of the three of the trials with SNR of 0 dB. This served the dual purpose of

132        ensuring that subjects could perform the task and familiarizing them with the talkers' faces and

133        voices.

134    Subjects subsequently completed 192 trials. At 25, 50, and 75 percent completion they were

135    given self-timed breaks with a minimum duration of 30 seconds. Each trial consisted of a video

136    of a talker saying a unique high probability sentence with the two background auditory streams

137    playing. No sentences were repeated. The trial began with a 2 second pause on the first frame of

138    the video providing the subject some time to process which voice to listen to for that trial. There

139    were 12 randomly interleaved conditions: three SNRs (0, −3, and −6 dB) and four visual

140    conditions (rotation of 0°, ±90°, or 180°, and a static upright image). After the video played,

141    subjects were instructed to type in what they heard with minimal spelling errors. They were also

142    informed that they would receive partial credit and to give a best guess if they were not certain.
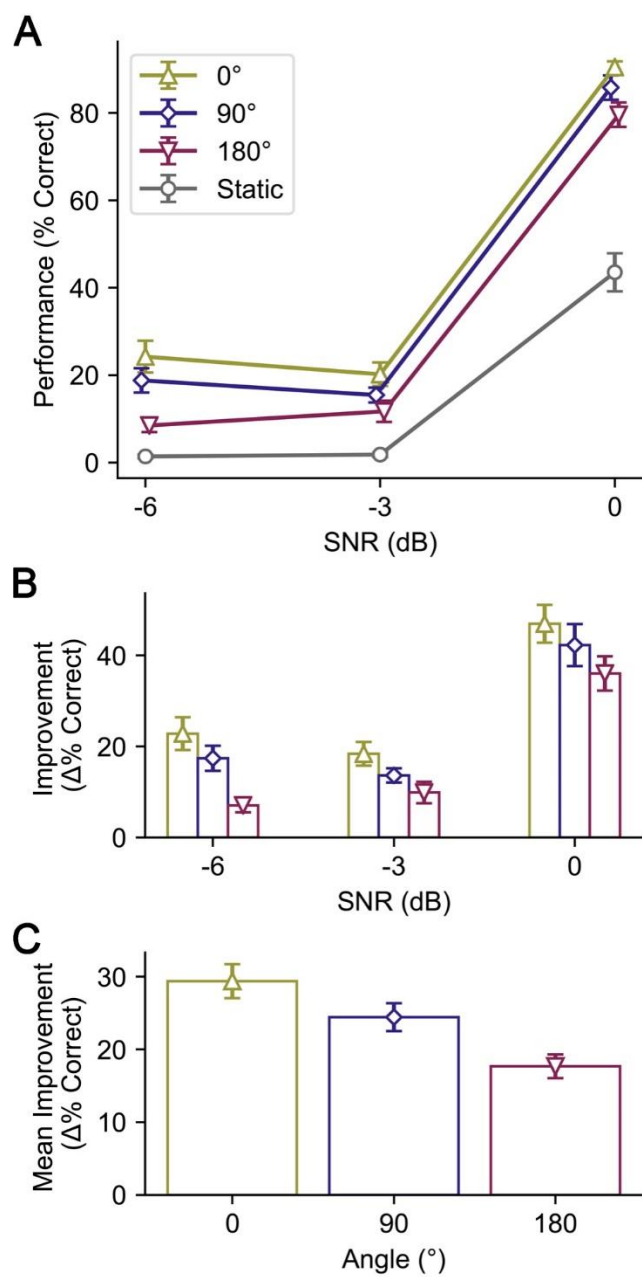
143  f.  *Scoring.* The sentences from the STeVi corpus had three to five predetermined keywords per

144    sentence. The percent of accurately entered keywords was hand scored based on Smayda et al's

145    (2016) criteria. Responses were considered correct if the spelling errors did not change the

146    meaning of the word or if the words were homophones.

147  g.  *Statistics.* We fit the data to a linear mixed effects model considering both SNR and angle to be

148    categorical variables (we did not expect a linear relationship, nor could we assume monotonicity,

149    particularly with regards to angle). We also considered interactions of SNR and angle. Each

150    subject was fit with an intercept.

151

152    **B. Results**

153    In the static condition, in which the face gives neither temporal nor linguistic information,

154    subjects had near chance performance in the speech in noise task at negative signal to noise ratios

155    (Figure 1A). Only at 0 dB SNR did subjects perform reliably above chance for the static face. Across

156    SNRs, subjects were able to significantly improve their performance when the face was moving by

157    7–47%, depending on the SNR and face rotation (Figure 1B). Across face rotation conditions these

158    gains were largest at 0 dB SNR and decreased as SNR worsened. Within each SNR, improvements

159    were largest for the upright face (0° rotation), and smallest for the inverted face (180° rotation).

160    Averaging across SNRs to specifically investigate the effect of disrupted linguistic cues, there were

161    significant differences in each subject's improvement depending on the rotation of the face. The

162    mixed effects model showed that 0 dB and −3 dB conditions were significantly different

163    ($p=1.15 \times 10^{-10}$), but there was not a difference between −3 dB and −6 dB. The 0° and 180° rotations

164    were significantly different($p=7.71 \times 10^{-3}$), as were the 90° and 180° rotations ($p=0.0119$). The

165    difference between 0° and 90° rotations approached, but did not reach, significance ($p=0.0876$).

166    There were no significant interaction terms.

167

11

168 Figure 2: (Color Online). A. Performance in the speech in noise task averaged across subjects.

169   Upward-facing triangle indicates the unrotated or upright face, diamond indicates rotation of

170   the face by 90° to the left or right, downward-facing triangle indicated the inverted face,

171   circle indicates a static image. B. The improvement in performance due to the moving face

172   calculated as the difference in each video condition and the respective static face for a given

173   SNR. C. Performance improvement due to temporal and linguistic cues averaged across

174   SNR conditions. All error bars show ±1 SEM.

175

176 **C. Summary**

177  In Experiment I we demonstrated that rotation of the head disrupted speech comprehension,

178 suggesting orientation specific processing of the face. Given the significant reduction in speech

179 processing between upright (0°) and inverted (180°) faces, we used these rotations to probe the

180 question of whether binding is affected by linguistic cues in Experiment II.

181 **III. EXPERIMENT II**

182 **A. Methods**

183  We engaged participants in a fundamental frequency modulation discrimination task. Listeners

184 were asked whether a target talker was modulated in pitch while ignoring masker talkers. They

185 simultaneously performed a visual detection catch trial task to ensure visual attention was

186 maintained throughout the experiment.

187 a. *Participants.* 23 participants (17 female, 6 male) ages 19−35 (mean 23.1) met the same criteria as

188  the participants from Experiment I. Six of the participants had also participated in Experiment I

189  and had similar performance to those who were naïve to the stimuli.

12

190    b.   *Stimuli.* Target and masker high context sentences were selected from the STeVi corpus. The

191      average duration of these sentences was 2.36 seconds with a standard deviation of 0.28 seconds.

192      Each trial included two of these sentences (one target and one masker sentence). These voice

193      pairings were evenly distributed across trials and always consisted of one male and one female

194      talker. The sentence pairings were randomly chosen, with each target sentence presented twice,

195      to have similar durations. All but nine sentence pairs had duration differences under 100 ms and

196      the maximum duration difference was 274 ms.

197      For some trials the audio from the videos was pitch modulated. Pitch modulations were 10

198      Hz cosine modulations with peak-to-peak amplitude of two semitones added to the stimulus'

199      natural pitch trajectory using Praat (Boesma and Weenick, n.d.). Videos were presented upright

200      or were inverted (rotated 180 degrees). The same audiobooks from Experiment I were played at

201      −6 dB SPL to provide additional interfering speech noise and make the task appropriately

202      challenging.

203    c.   *Procedure.* Subjects were given three chances to pass a training module in which they had to detect

204      pitch modulation when no maskers were present. As in Experiment I, this allowed subjects to

205      learn the identities of the talkers. They were given ten practice trials and ten testing trials for

206      which they had to achieve 70% accuracy to pass. The statistics of modulation were the same for

207      both the training and the main experiment. Subjects were then shown two example trials to

208      familiarize them with the complexity of the stimulus. They were instructed to look at the faces

209      on the screen and perform two tasks simultaneously: the main pitch modulation discrimination

210      task, and a catch trial visual detection task. Breaks were offered as in Experiment I.

211    d.  *Main Task.* For each trial there was a target talker and masker talker. At the beginning of the trial,

212       the subject saw an image of the target talker for 1.5 seconds, indicating which to listen to during

213       the trial. After the image was presented, the video and four auditory streams were played: the

214       target talker, the masker talker, and the two interfering audiobook background streams. Subjects

215       reported whether the voice of the target talker contained the modulation by pressing a button.
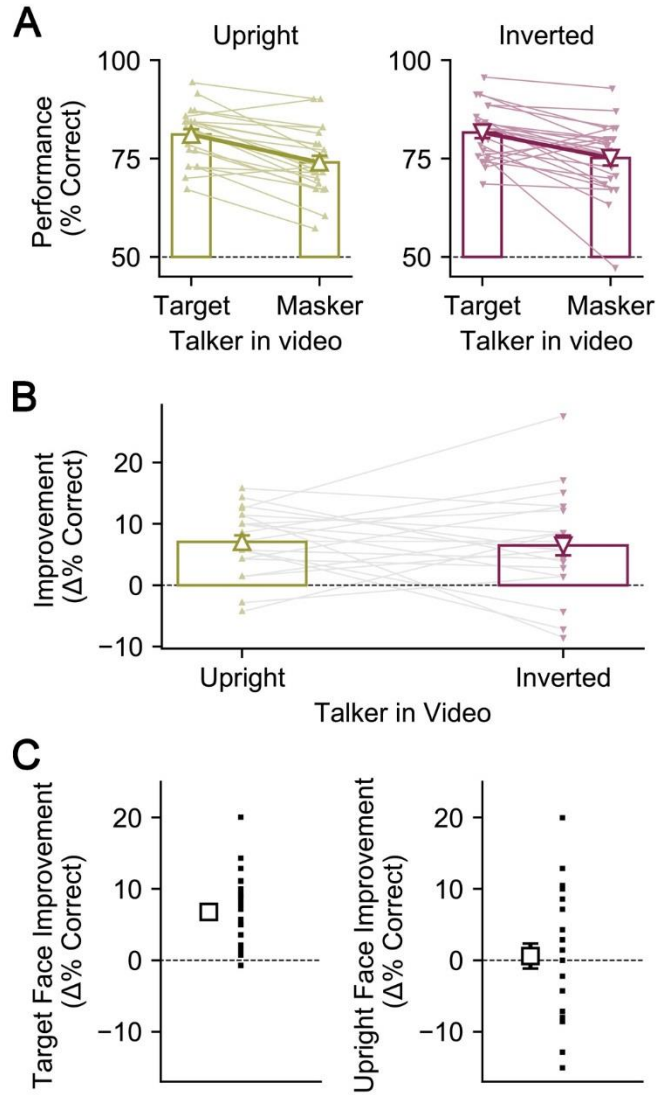
216          This task consisted of 16 conditions (2 masker/target video × 2 face rotations × 4

217       masker/target modulations). The video of a talker either matched the target or the masker

218       auditory stream, the face of the talker was upright or inverted, and both, neither, only the target,

219       or only the masker auditory streams were modulated. The subjects heard each of the 200 high

220       probability sentences twice as a target, for a total of 400 trials. There were 35 trials for each

221       condition in which only the masker or target were modulated (70% of trials; the conditions of

222       interest) and 15 trials for remaining conditions which were included so the subject could not

223       infer the target modulation based on the masker modulation.

224    e.  *Catch Trials.* Of the 120 trails for which either both or neither sentences were modulated, 36

225       random trials had a small pink translucent dot over the mouth of the talker. Subjects were

226       instructed to press 3 when they saw this dot and to not respond for pitch modulation. These

227       catch trials ensured that subjects were looking at the talker's face. Subjects were informed each

228       time they failed to detect the dot. The criterion to be included in the analysis was detection of

229       more than 80% of the dots. All subjects achieved this, so all data were included in the analysis.

230  **B.  Results**

231      All subjects were able to perform the pitch discrimination task well above chance (Figure 2A).

232      There was a significant improvement in both the upright (paired t-test, t=6.61, p=1.20x10$^{-6}$) and

233      inverted (paired t-test, t=3.91, p=7.49x10$^{-4}$) conditions when subjects viewed the target face relative

234      to when they viewed the masker face (Figure 2B). However, there was no difference in the benefit

235      due to the target face between the upright and inverted condition, and therefore no benefit of the

236      upright face (Figure 2C). Averaging across face rotation conditions, subjects experienced a

237      significant benefit (paired t-test, t=6.61, p=1.20x10$^{-6}$) of approximately 7% when the face of the

238      video matched the target rather than the masker.

15

240

241    Figure 3: (Color Online). A. Performance in the pitch modulation discrimination task. Small solid

242        markers show individual subjects, larger open markers show the average across subjects. B.

243        The improvement in performance due to the video of the target talker's face calculated as

244        the difference between the target visual condition and masker visual condition for a given

245        face rotation. C. (Left) Net improvement due to the target talker's face (difference between

246        target and masker conditions averaged across subjects and upright/inverted conditions).

247        (Right) Net improvement due to the upright face (difference between upright and inverted

248        conditions averaged across subjects and target/masker conditions). All error bars show ±1

249        SEM.

250

251    **C.  Summary**

252        In Experiment II the rotation of the face did not significantly influence binding. Nonetheless

253    this paradigm showed a strong replication of previous findings that temporal coherence induces

254    binding (Atilgan et al., 2018; Maddox et al., 2015). The importance of temporal coherence for

255    binding has not previously been established for speech.

256    **IV.    DISCUSSION**

257        Together these experiments suggest that visual linguistic cues and audio-visual binding

258    contribute independently to processing multimodal speech in noise. Experiment I addressed the

259    benefit of visual linguistic cues in speech comprehension and Experiment II investigated audio-

260    visual binding, ultimately showing that visual linguistic cues do not enhance the listener's object-

261    based attention to the target talker.

262      In Experiment I we demonstrated that both temporal and linguistic cues are important for

263      speech comprehension in a speech in noise task. There was a significant improvement in

264      performance when the face was moving relative to the static image. This was true of all face rotation

265      conditions. Because temporal cues were preserved across rotation conditions, we consider some

266      portion of the video performance improvement to be due to temporal cues. However, as the

267      rotation of the face increased in magnitude, the benefit of the video decreased, suggesting that some

268      of the benefit in each condition is due to processing of linguistic cue. Though linguistic cues are

269      present even in the rotated faces, the processing of this information seems to be impaired by

270      rotating the face. Interestingly, performance drops with the magnitude of the rotation even though

271      the 0° and 180° rotation are more geometrically similar due to their vertical symmetry than the 0°

272      and 90° rotations. There are two possible explanations for this: the subject can partially compensate

273      for the face's rotation when processing visemes or that the subject has more prior experience with

274      90° rotated faces than with 180° rotated faces. While the latter is likely true, we do not believe it is a

275      compelling explanation for our results. A vast majority of conversations are held with upright faces,

276      and situations in which we are speaking to someone at a 90° rotation are minimal (e.g. talking to

277      someone while reclined). Therefore, it seems more likely that subjects are "un-rotating" the face

278      where possible to get some benefit from linguistic cues, and this is easier for them to do with 90°

279      rotation than 180°.

280    In Experiment II we show that fundamental frequency modulation discrimination is improved

281    when listeners can see the video of the target talker rather than a masker talker regardless of the

282    orientation of the face. Structurally our experiment was very similar to previous work that tested for

283    binding by engaging listeners in simultaneous auditory discrimination and visual detection tasks

284    (Maddox et al., 2015). Importantly, the tasks rely on the tracking of an orthogonal perceptual feature

285    (pitch), one that is independently changing to the feature that is coherently modulated. If the listener

286    binds the auditory and visual streams based on their temporal coherence, their brain will form a

287    perceptual object. By allocating object-based attention, all features of the object, including the

288    orthogonal feature will be enhanced, leading to better performance. The performance improvement

289    is not explained by late-stage integration since the visual stream provides no information about the

290    orthogonal auditory features.

291    We improved upon Maddox et al's original task (2015) by using more natural stimuli. In this case

292    the listeners had to simultaneously determine whether there was a pitch modulation in the target

293    talker or a pink disk on the mouth in the video while ignoring the masker talker. We used real

294    speech as the stimuli and pitch as the orthogonal feature, which gave ecological relevance to the

295    task. Processing of pitch modulations are important in natural environments due to prosodic

296    information that is in part carried by the pitch of a talker. This prosodic information not only

297    provides emotional context but also is important for parsing full sentences (Stirling, 1996; Warren et

298    al., 1995). Binding of audiovisual speech could improve our perception of not only what the talker is

299    saying, but how they are saying it. Binding can explain an improvement in performance when the

300    video matches the target. The listener can allocate object-based attention to the target talker and

301    improve their discrimination of pitch modulation because detection of visual events will not divide

302    their attention.

303      There is a consistent improvement in processing orthogonal stimulus features when the listener

304    can see the target video, which can be explained by binding. However, the benefit is not modulated

305    by rotating the face, suggesting that temporal coherence is the cue that underpins binding in this

306    experiment. Using real speech, our results confirm the finding that temporal coherence drives

307    audiovisual binding, which had been previously established for stimuli with speech-like dynamics

308    (Maddox et al., 2015).

309      We did not find an effect of face rotation on performance in the pitch discrimination task. There

310    are a few possible explanations for this. Temporal coherence may be sufficient to induce binding,

311    and a possible contribution of linguistic cues would be overshadowed by the influence of temporally

312    coherent cues. Alternatively, a contribution of linguistic cues to strengthen binding, if such a thing is

313    possible, may have been too small to be measured behaviorally. If linguistic cues truly do not

314    influence binding, the hierarchical processing of language therefore suggests an explanation for

315    binding occurring independent of face rotation. While low level spectral features are well

316    represented in A1, articulatory features are not represented until the superior temporal gyrus (STG)

317    (Ding et al., 2016; Mesgarani et al., 2014). A study involving ferrets performing a multisensory task

318    found neural evidence of binding in primary auditory cortex (A1) (Atilgan et al., 2018), whereas

319    traditional Bayesian or late-stage integration is thought to occur at higher processing areas in the

320    intraparietal sulcus (Rohe and Noppeney, 2015, 2016). Such findings support the notion that binding

321    and late-stage integration are fundamentally different processes that rely on different types of

322    sensory information. The extent of the visual and auditory information available at such early

323    processing areas to create binding is uncertain, particularly given the unknown origins of the visual

324    connections responsible for visual-dependent auditory activity in A1. In order for linguistic cues to

325    contribute to binding the brain would need to combine feedback from STG carrying auditory

326    articulatory information with viseme information from visual areas.

327    Together these experiments show that face-rotation and therefore disruption of linguistic cues

328    hinders audiovisual speech comprehension, but not detection of orthogonal pitch features. Even if

329    linguistic cues do play a role in binding, their behavioral benefit seems to be superseded by temporal

330    coherence. Therefore, the benefit of visual linguistic cues to speech understanding is likely due to

331    late integration in which visemes can bias the listener towards the correct phoneme perception at

332    higher processing stages. Binding, then, may be specific to very low-level physical correspondences,

333    a hypothesis on which future experiments will shed more light.

334    **V.      CONCLUSION**

335    We demonstrated the importance of both temporal and linguistic visual cues for audiovisual

336    speech in noise comprehension in an ecologically relevant task. We also showed that audiovisual

337    temporal coherence, but not linguistic congruence, improved performance in a frequency

338    modulation discrimination task, consistent with the existence of audiovisual binding. It thus appears

339    that multisensory linguistic cues, an example of learned physical correspondence, are integrated at

340    the perceptual decision-making stage rather than early integration. Practically, our results suggest that

341    visemes can benefit listeners in noisy environments by biasing the listener towards perceiving the

342    correct sentence, but they do not aid listeners in detecting other aspects of the talker's speech.

343    **ACKNOWLEDGMENTS**

347    **REFERENCES**

348    Arnold, P., and Hill, F. (**2001**). "Bisensory augmentation: A speechreading advantage when speech is

349        clearly audible and intact," Br. J. Psychol., **92**, 339–355. doi:10.1348/000712601162220

350    Atilgan, H., Town, S. M., Wood, K. C., Jones, G. P., Maddox, R. K., Lee, A. K. C., and Bizley, J. K.

351         (**2018**). "Integration of Visual Information in Auditory Cortex Promotes Auditory Scene

352         Analysis through Multisensory Binding," Neuron, **97**, 640-655.e4.

353         doi:10.1016/j.neuron.2017.12.034

354    Bizley, J. K., Maddox, R. K., and Lee, A. K. C. (**2016**). "Defining Auditory-Visual Objects:

355         Behavioral Tests and Physiological Mechanisms," Trends Neurosci., **39**, 74–85.

356         doi:10.1016/j.tins.2015.12.007

357    Blaser, E., Pylyshyn, Z. W., and Holcombe, A. O. (**2000**). "Tracking an object through feature

358         space," Nature, **408**, 196-.

359    Boesma, P., and Weenick, D. (**n.d.**). *Praat: doing phonetics by computer,.*

360    Ding, N., Melloni, L., Zhang, H., Tian, X., and Poeppel, D. (**2016**). "Cortical tracking of hierarchical

361         linguistic structures in connected speech," Nat. Neurosci., **19**, 158–164. doi:10.1038/nn.4186

362    Grant, K. W., and Bernstein, J. G. W. (**2019**). "Toward a Model of Auditory-Visual Speech

363         Intelligibility," In A. K. C. Lee, M. T. Wallace, A. B. Coffin, A. N. Popper, and R. R. Fay

364         (Eds.), Multisensory Process. Audit. Perspect., Springer Handbook of Auditory Research,

365         Springer International Publishing, Cham, pp. 33–57. doi:10.1007/978-3-030-10461-0_3

366    Körding, K. P., Beierholm, U., Ma, W. J., Quartz, S., Tenenbaum, J. B., and Shams, L. (**2007**).

367         "Causal Inference in Multisensory Perception," PLOS ONE, **2**, e943.

368         doi:10.1371/journal.pone.0000943

369    Lee, A. K. C., Maddox, R. K., and Bizley, J. K. (**2019**). "An Object-Based Interpretation of

370         Audiovisual Processing," In A. K. C. Lee, M. T. Wallace, A. B. Coffin, A. N. Popper, and R.

371         R. Fay (Eds.), Multisensory Process. Audit. Perspect., Springer Handbook of Auditory

372         Research, Springer International Publishing, Cham, pp. 59–83. doi:10.1007/978-3-030-

373         10461-0_4

374    L'Engle, M. (**2006**). *A Wrinkle in Time*, Findaway World Llc, 238 pages.

375    Maddox, R. K., Atilgan, H., Bizley, J. K., and Lee, A. K. (**2015**). "Auditory selective attention is

376        enhanced by a task-irrelevant temporally coherent visual stimulus in human listeners," eLife,

377        , doi: 10.7554/eLife.04995. doi:10.7554/eLife.04995

378    Massaro, D. W., and Cohen, M. M. (**1996**). "Perceiving speech from inverted faces," Percept.

379        Psychophys., **58**, 1047–1065. doi:10.3758/BF03206832

380    Mcgurk, H., and Macdonald, J. (**1976**). "Hearing lips and seeing voices," Nature, **264**, 746–748.

381        doi:10.1038/264746a0

382    Mesgarani, N., Cheung, C., Johnson, K., and Chang, E. F. (**2014**). "Phonetic Feature Encoding in

383        Human Superior Temporal Gyrus," Science, **343**, 1006–1010. doi:10.1126/science.1245994

384    Reisberg, D., McLean, J., and Goldfield, A. (**1987**). "Easy to hear but hard to understand:  A lip-

385        reading advantage with intact auditory stimuli," Hear. Eye Psychol. Lip-Read., Lawrence

386        Erlbaum Associates, Inc, Hillsdale, NJ, US, pp. 97–113.

387    Rohe, T., and Noppeney, U. (**2015**). "Cortical Hierarchies Perform Bayesian Causal Inference in

388        Multisensory Perception," (C. Kayser, Ed.) PLOS Biol., **13**, e1002073.

389        doi:10.1371/journal.pbio.1002073

390    Rohe, T., and Noppeney, U. (**2016**). "Distinct Computational Principles Govern Multisensory

391        Integration in Primary Sensory and Association Cortices," Curr. Biol., **26**, 509–514.

392        doi:10.1016/j.cub.2015.12.056

393    Scott, M. (**2008**). *The Alchemyst: The Secrets of the Immortal Nicholas Flamel*, Delacorte Press, 402 pages.

394    Shinn-Cunningham, B. G. (**2008**). "Object-based auditory and visual attention," Trends Cogn. Sci.,

395        **12**, 182–186. doi:10.1016/j.tics.2008.02.003

396    Smayda, K. E., Engen, K. J. V., Maddox, W. T., and Chandrasekaran, B. (**2016**). "Audio-Visual and

397        Meaningful Semantic Context Enhancements in Older and Younger Adults," PLOS ONE,

398        **11**, e0152773. doi:10.1371/journal.pone.0152773

399    Stevenson, R. A., Zemtsov, R. K., and Wallace, M. T. (**2012**). "Individual Differences in the

400        Multisensory Temporal Binding Window Predict Susceptibility to Audiovisual Illusions," J.

401        Exp. Psychol. Hum. Percept. Perform., **38**, 1517–1529. doi:10.1037/a0027339

402    *STeVi Speech Test Video Corpus* (**n.d.**). Sensimetrics. Retrieved from

403        https://www.sens.com/products/stevi-speech-test-video-corpus/

404    Stirling, L. (**1996**). "Does Prosody Support or Direct Sentence Processing?," Lang. Cogn. Process.,

405        **11**, 193–212. doi:10.1080/016909696387268

406    Ujiie, Y., Asai, T., and Wakabayashi, A. (**2018**). "Individual differences and the effect of face

407        configuration information in the McGurk effect," Exp. Brain Res., **236**, 973–984.

408        doi:10.1007/s00221-018-5188-4

409    Van Engen, K. J., Phelps, J. E. B., Smiljanic, R., and Chandrasekaran, B. (**2014**). "Enhancing speech

410        intelligibility: interactions among context, modality, speech style, and masker," J. Speech

411        Lang. Hear. Res. JSLHR, **57**, 1908–1918. doi:10.1044/JSLHR-H-13-0076

412    Warren, P., Grabe, E., and Nolan, F. (**1995**). "Prosody, phonology and parsing in closure

413        ambiguities," Lang. Cogn. Process., **10**, 457–486. doi:10.1080/01690969508407112

414