

## Acoustic facilitation of object movement detection during self-motion

Journal:	<i>Proceedings B</i>
Manuscript ID:	Draft
Article Type:	Research
Date Submitted by the Author:	n/a
Complete List of Authors:	Calabro, Finnegan; Boston University, Brain & Vision Research Lab, Department of Biomedical Engineering Soto-Faraco, Salvador; Universitat Pompeu Fabra, ICREA & Dept. de Tecnologies de la Informació i les Comunicacions Vaina, Lucia; Boston University, Brain & Vision Research Lab, Department of Biomedical Engineering
Subject:	Behaviour < BIOLOGY, Neuroscience < BIOLOGY
Keywords:	flow parsing, visual search, multisensory perception, auditory motion, visual motion
Proceedings B category:	Behaviour


 SCHOLARONE™  
Manuscripts

1 **Acoustic facilitation of object movement detection during self-motion**

2 Calabro FJ<sup>1</sup>, Soto-Faraco S<sup>2</sup>, Vaina, LM<sup>1,3</sup>

3

4 <sup>1</sup>Brain and Vision Research Lab, Department of Biomedical Engineering, Boston, MA,  
5 USA 02215

6 <sup>2</sup>ICREA & Dept. de Tecnologies de la Informació i les Comunicacions, Universitat  
7 Pompeu Fabra, Barcelona, Spain

8 <sup>3</sup>Departments of Neurology and Radiology, Massachusetts General Hospital, Harvard  
9 Medical School, Boston, MA, USA, 02215, [vaina@bu.edu](mailto:vaina@bu.edu)

10

For Review Only

11 **Abstract**

12 In humans, as well as most animal species, perception of object motion is critical to  
13 successful interaction with the surrounding environment. Yet, as the observer moves, the  
14 retinal projections of the various motion components add to each other and extracting  
15 accurate object motion becomes computationally challenging. Recent psychophysical  
16 studies have demonstrated that observers use a flow parsing mechanism to estimate and  
17 subtract self-motion from the optic flow field. We investigated whether concurrent  
18 acoustic cues for motion can facilitate visual flow parsing, thereby enhancing the  
19 detection of moving objects during simulated self-motion. Participants identified an  
20 object (the target) that moved either forward or backward within a visual scene  
21 containing nine identical textured objects simulating forward observer translation. We  
22 found that spatially co-localized, directionally congruent, moving auditory stimuli  
23 enhanced object motion detection. Interestingly, subjects who performed poorly on the  
24 visual-only task benefited more from the addition of moving auditory stimuli. When  
25 auditory stimuli were not co-localized to the visual target, improvements in detection  
26 rates were weak. Taken together, these results suggest that parsing object motion from  
27 self-motion induced optic flow can operate on multi-sensory object representations.

28  
29 *Keywords:* flow parsing, visual search, multisensory perception, visual motion, auditory  
30 motion

## 31 **Introduction**

32 For a stationary observer, an object moving within an otherwise still scene is uniquely  
33 identified by motion and can be effortlessly detected no matter how many elements the  
34 scene contains [1]. It has been suggested that motion responsive mechanisms filter out the  
35 static objects thus making detection of the unique moving elements trivial [2]. If the  
36 background objects can be grouped into a rigid surface defined by disparity [3] or  
37 common motion [4] then a single moving object will also pop-out. However, from the  
38 point of view of the visual system, static or rigid backgrounds are only an exceptional  
39 case given that an observer's translation and head motion produce a complex movement  
40 pattern of the visual field, or optic flow [5, 6]. This adds remarkable complexity when  
41 trying to distinguish object motion from the dynamic scene given that all objects move in  
42 terms of their retinal projections. Yet, perception of object motion during self-movement  
43 in humans is both accurate as well as critical to successful interaction with the  
44 environment. It has been proposed that object motion can be parsed out from the optic  
45 flow created by self motion, thus allowing a moving observer to detect a moving object.  
46 Rushton and colleagues [7-9] have suggested a flow-parsing mechanism that uses the  
47 brain's sensitivity to optic flow to separate retinal motion signals into those components  
48 due to observer movement and those due to the movement of objects in the scene.

49 Previous studies addressing flow parsing have concentrated on the visual modality  
50 alone. Although vision is dominant in our perception of motion, natural environments  
51 frequently provide extra-visual cues to motion, such as the sound of a car down the street  
52 quickly approaching to, or moving away from, us. The question addressed here is  
53 therefore whether extra-visual (in this case, acoustic) information can complement optic

54 flow parsing, and hence facilitate the extraction of visual motion from dynamic visual  
55 scenes during observer movement. The benefits of congruent, cross-modal information  
56 are well known, especially the enhancement of responses to a stimulus in situations when  
57 the signal from a single modality is weak [10-12] or when processing within one sensory  
58 system is impaired by brain damage [13]. In the particular case of motion, strong  
59 synergies between different sensory modalities have been described in several recent  
60 studies. For example, in horizontal motion, directional incongruence between visual and  
61 auditory signals can lead to strong illusions in the direction of the sound (e.g., [14-16] see  
62 [17] for a review), whereas directional congruence can lead to improved detection  
63 performance (e.g., [18-20], though [19] suggests the improvement may be statistical,  
64 rather than due to bimodal enhancement). Speed of motion is subject to similar  
65 phenomena, given that sounds will appear to move faster (or slower) depending on the  
66 velocity of concurrent visual stimuli [21].

67         Directly relevant to the present question, several previous studies showed that  
68 cross-modal directional congruency effects can be observed in motion along the depth  
69 plane. For example, auditory looming has been shown to speed up the detection of  
70 looming visual signals [22], and an appropriately timed static sounds can drive the  
71 perceived direction of an ambiguous visual apparent motion stimulus [23]. In studies  
72 using motion after-effects in the depth plane [24, 25], adaptation to looming (or receding)  
73 visual stimuli produced an after-effect in the reverse direction for subsequently presented  
74 sounds. When using directionally incongruent audiovisual adaptors, the after-effect is  
75 consistent with the direction of the visual adaptor. The phenomenology of these findings  
76 suggests that the interaction between visual and auditory motion signals can express at

77 rather early levels of processing. Indeed, recent studies using fMRI have revealed that  
78 cross-modal motion congruency effects are reflected in a complex network of brain  
79 structures including unisensory motion processing areas as well as areas of multisensory  
80 convergence [26, 27]. In particular, illusory reversals of sound direction (induced by  
81 directional incongruence between auditory and visual motion) were correlated with a  
82 deactivation of auditory motion areas (the Auditory Motion Complex, AMC) and an  
83 enhancement of activity in the cortical areas responsive to visual motion. Furthermore, in  
84 the same study it was shown that, just prior to trials leading to illusory sound motion  
85 percepts, activity in the ventral intraparietal area (VIP, an area of multisensory  
86 convergence that contains spatial representations) was stronger than in identical trials that  
87 resulted in veridical perception of sound direction [26].

88 We used a visual search paradigm [28-32] that has been previously used to test  
89 optic-flow parsing [1-3, 7, 33-35]. Several recent findings attest to the potential of cross-  
90 modal enhancement by sounds to improve visual selective attention in search tasks [36,  
91 37]. For example, Van der Burg et al. (2008) showed that sounds temporally coincident  
92 with an irrelevant change in visual targets dramatically improved search times in a  
93 difficult search task. In fact, a difficult visual search task which led to serial search  
94 patterns in the absence of sounds, reflected nearly flat search slopes when a sound was  
95 synchronized target changes. Interestingly, when the sound was paired with a visual  
96 distractor instead, the search became more difficult. These demonstrations, together with  
97 the strong cross-modal synergies in motion processing described above, highlight the  
98 possibility that acoustic motion could help parse out object motion from optic flow from  
99 in dynamic visual displays.

100 Note, however, that none of the previous studies addressing cross-modal  
101 enhancement in visual search have, to our knowledge, involved dynamic scenes.  
102 Moreover, paradigms where perceptual load is high (i.e., when the matching between  
103 sound and visual events must be extracted from complex, dynamically changing events)  
104 have typically failed to demonstrate cross-modal enhancement in search tasks [38, 39].

105 It is therefore uncertain whether the visual motion processes leading to parsing  
106 out object motion from optic flow produced by the observer's movement can benefit  
107 from cross-modal synergies. Here we addressed this question empirically. We compared  
108 performance on a task of object movement detection during self-motion when paired with  
109 a static or moving auditory cue to determine whether cross-modal motion congruency  
110 enhances visual selection. Our results show that while auditory stimuli not co-localized  
111 with the visual target impart only a small benefit to detection rates, the presentation of a  
112 moving, co-localized auditory cue provides a significant gain.

### 113 **Methods**

#### 114 **Subjects**

115 All participants (n=18, 8 male, age range 19-29, mean 22) performed the visual  
116 task, and each was tested with either the non-co-localized (n=10) or co-localized (n=10)  
117 auditory condition. Two of the participants including FC (an author) performed both  
118 auditory conditions, and all (except for FC) were naïve to the purposes of the experiment.  
119 All procedures were approved by the Boston University Institutional Review Board for  
120 research with human subjects, and informed consent was obtained from each participant.

**121 Apparatus**

122 Participants viewed the visual display from a distance of 60 cm, with head  
123 position fixed by a chin and forehead rest. Stimuli were displayed on a 23" Apple  
124 Cinema Display and were generated in Matlab using Psychophysical Toolbox [40, 41]  
125 and OpenGL libraries. Suprathreshold auditory cues were presented with Bose QC-1  
126 QuietComfort acoustic noise canceling headphones. We used a Minolta LS-100 for  
127 monitor luminance calibration, and a Scantek Castle GA-824 Smart Sensor SLM for  
128 acoustic calibration.

**129 Stimulus**

130 Participants viewed nine textured spheres distributed within a simulated virtual  
131 environment of size 25x25x60 cm (Figure 1A), projected onto an Apple Cinema Display.  
132 Stimuli were viewed binocularly, but contained no stereo cues, such that visual motion in  
133 depth was determined only looming motion. To avoid overlapping spheres, the viewable  
134 area was divided into nine equally sized wedges in the frontoparallel plane, and one  
135 object was placed into each wedge with a randomly chosen eccentricity. Objects were  
136 located randomly in simulated depth between 25 and 35 cm, and had a mean diameter of  
137  $1.5^\circ$  at the start of the stimulus, with a mean luminance of  $28 \text{ cd/m}^2$  on a background of  
138 luminance  $0.3 \text{ cd/m}^2$ . A red fixation mark was placed at the center of the display and  
139 subjects were instructed to maintain fixation throughout the testing block.

140 Forward observer motion was simulated towards the fixation mark for 1 sec.  
141 Except where noted, the observer motion was a forward translation of 3 cm/sec (thus  
142 inducing a corresponding expansion of objects that were stationary within the scene).  
143 One of the nine objects, the target, was assigned an independent motion vector, moving

144 either forward or backward at 2, 4, 6, or 8 cm/sec with respect to the rest of the scene  
145 (Figure 1B). The target's visible motion was the sum of its own motion vector and the  
146 induced motion caused by the simulated translation of the observer. At the end of the  
147 motion the screen was cleared for 250 ms before all objects reappeared at their final  
148 locations, but projected into a single depth plane so that all were a constant size. Four  
149 objects (the target and 3 other randomly selected spheres) were shown with labels  
150 (marked with numerals, 1-4) and observers were asked to report which of the four was  
151 the one that had been moving within the scene (and not solely because of the observer  
152 translation). Since the labels appeared only after the end of the trial, subjects had to  
153 monitor all nine objects, although their decision was a 4 alternative forced choice (4AFC)  
154 task.

155 In separate conditions, the stimulus was presented visually only or with either a  
156 co-localized or non-co-localized (central) auditory cue. The auditory cue was a pure tone  
157 of frequency 300 Hz that, in 75% of trials (*auditory-moving* trials), was simulated (via a  
158 change in amplitude) to move within the scene in the same direction as the target  
159 (forward or backward), and in the remaining 25% of trials was presented at constant  
160 amplitude throughout the trial (*auditory-static* trials). The change in amplitude was  
161 modeled as a sound source at an initial distance of 4.1 meters (69 dB SPL) moving  
162 towards or away from the observer at 3.5 m/s (resulting in a change of approximately 10  
163 dB SPL, Figure 1C). Sound attenuation as a function of distance was approximated for  
164 the testing room by measuring sound levels at various distances from a constant sound  
165 source. A least-squares fit was applied to determine the relationship of sound amplitude  
166 to distance. In both the static and moving (whether approaching or receding) auditory

167 conditions, the auditory cue started at the same amplitude so that the initial volume did  
168 not indicate whether the auditory cue would move, or in what direction. In all auditory  
169 conditions, the auditory stimulus was enveloped with 30 ms ramps to avoid clicks due to  
170 a sharp onset or offset. Participants were screened to ensure they could discriminate the  
171 direction of the auditory motion.

172 In the non-co-localized auditory condition, the auditory cue appeared to arise  
173 directly in front of the observer (it was presented with equal amplitude to both ears). In  
174 the co-localized auditory condition, the inter-aural intensity difference (IID) was adjusted  
175 to match the horizontal eccentricity of the target object. For both the non-co-localized  
176 and co-localized auditory cues, we used auditory-moving and auditory-static conditions  
177 to distinguish effects due to localizing the target, effects due to congruent auditory  
178 motion and effects which require both spatially co-localized and congruent-motion  
179 auditory cues.

180 ----- Figure 1 near here -----

## 181 **Results**

### 182 **Detection of object movement during self-motion**

183 Figure 2 shows the results from all 18 subjects on the visual-only condition. As  
184 expected, performance depended on the speed of the target object, with faster speeds (6  
185 and 8 cm/sec) being detected above 80% correct. Performance was above chance  
186 (chance=25%) for all speeds. Approaching objects (those moving towards the observer  
187 within the simulated environment) were easier to detect than receding ones, as  
188 demonstrated by the increased performance for positive speeds relative to negative

189 speeds: a 2-way ANOVA showed significant effects of target speed ( $F_{(3, 152)}=143.5$ ,  
190  $p<0.001$ ) and direction ( $F_{(1, 152)}=47.5$ ,  $p<0.001$ ).

191 ----- Figure 2 near here -----

## 192 **How is object motion detected?**

193 Object motion detection in the visual search task may be accomplished by flow-  
194 parsing, as suggested by Rushton and colleagues [7-9], in which self-motion is estimated  
195 from background optic flow and parsed-out from the scene. Alternatively, to resolve this  
196 task, participants may use the object's motion relative to the observer (i.e., retinal  
197 motion), for example detecting an object with a high perceived speed, or an object which  
198 appears nearly stationary among moving objects (as in [42]). To determine which of these  
199 mechanisms was most likely used in our experiment, ten participants performed an  
200 additional visual-only condition in which the speed of the simulated observer motion was  
201 increased to 5 cm/sec. If observers used relative motion cues, this should have resulted in  
202 worse performance for the -6 cm/sec target speed (where the target's retinal speed  
203 decreased with the faster observer motion), and better performance for objects with  
204 positive (2, 4, 6 cm/sec) velocities (in which the retinal speed increased with the faster  
205 observer motion). If, on the other hand, subjects used a flow parsing mechanism,  
206 performance levels should have been independent of the self-motion speed (which is  
207 parsed out), as long as self-motion was easily detected, as was the case in both observer  
208 speed conditions (3, 5 cm/sec).

209 Figure 3 shows the results for an observer speed of 5 cm/s compared to results  
210 from the same ten subjects when the observer speed was 3 cm/sec. A two-way ANOVA  
211 showed a significant effect of target speed ( $F_{(5, 113)}=84.4$ ,  $p<0.001$ ), thus reproducing the

212 result of the main visual-only experiment, but no effect of observer speed ( $F_{(1, 113)}=0.12$ ,  
213  $p>0.7$ ). We further tested the two predictions of the retinal motion hypothesis separately:  
214 (1) a decrease in performance at -6 cm/sec due to lower retinal speed at the higher  
215 observer speed, and (2) an increase in performance for positive object speeds due to the  
216 increase in retinal speed at the higher observer speed. A t-test considering only data from  
217 the -6 cm/sec object speed (prediction 1), showed no difference with changes in observer  
218 motion ( $t(9)=-0.11$ ,  $p=0.91$ ), and a two-way ANOVA restricted to positive target speeds  
219 (prediction 2) similarly showed no significant effect of observer speed ( $F_{(1, 56)}=0.13$ ,  
220  $p>0.7$ ). Furthermore, a two-one-sided t-test (TOST) [43, 44] for equivalence showed that  
221 performance for the two observer speed conditions across subjects and object speed was  
222 statistically equivalent at  $p<0.05$  within a tolerance of 2.5%. Since a change of 2 cm/sec  
223 caused on average a 21% change in performance when applied to the object speed,  
224 equivalent performance within a 2.5% tolerance when the 2 cm/sec speed difference was  
225 applied to the observer speed indicates that the difference in retinal motion speeds cannot  
226 account for performance on the visual task. Taken together, these results suggest it was  
227 unlikely that observers solved the task by using only retinal motion cues.

228 ----- Figure 3 near here -----

### 229 **Do auditory cues facilitate detection of object movement during self-motion?**

230 To determine whether auditory motion cues can facilitate the detection of object  
231 movement, we considered two conditions in which a moving auditory cue was presented  
232 with motion direction congruent to that of the visual target. First, we tested whether the  
233 detection of object movement during self-motion is facilitated by the presentation of a  
234 synchronous, but spatially non-co-localized, auditory cue (perceptually located at the

235 center of the display). Second, we tested whether facilitation depends on the spatial co-  
236 localization of the visual and auditory motions (with an IID matching the horizontal  
237 eccentricity of the visual target). In both cases, performance was compared to that of  
238 static auditory cues (non-co-localized and co-localized, respectively).

### 239 *Non-co-localized auditory stimulus*

240 Figure 4 shows the performance of 10 subjects on the moving object detection  
241 task in the presence of a non-co-localized auditory cue (localized to the center of the  
242 screen). A two-way ANOVA showed a small, non-significant increase in performance  
243 (3.2% mean improvement) for *auditory-moving* trials ( $F_{(1, 144)}=3.39$   $p=0.06$ ) as compared  
244 to *auditory-static* trials. There was a significant main effect of target speed ( $F_{(7, 144)}=81.9$ ,  
245  $p<0.001$ ), but no significant interaction between auditory condition (static vs. moving)  
246 and target speed ( $F_{(7, 144)}=0.64$ ,  $p>0.7$ ). These results suggest that the presentation of a  
247 synchronous auditory cue that is not spatially co-localized with the target produced only a  
248 very modest improvement in the detection of a moving object.

249 ----- Figure 4 near here -----

250 An analysis of reaction times in trials with correct responses showed that both  
251 auditory-static and auditory-moving trials resulted in faster response times than the visual  
252 only condition in the same subjects, by 43 ms ( $F_{(1, 57)}=10.12$ ,  $p=0.002$ ) and 41 ms ( $F_{(1, 57)}=26.59$ ,  
253  $p<0.001$ ) respectively. However, there was no significant difference between  
254 the auditory static and auditory motion conditions ( $F_{(1, 59)}=0.14$ ,  $p>0.7$ ). Therefore, the  
255 use of a non-co-localized auditory motion cue contributed neither a statistically  
256 significant increase in performance, nor decrease in response time.

257 *Co-localized auditory stimulus*

258 To test the effect of spatial co-localization on auditory facilitation in our task, 10  
259 participants performed a version of the task in which the IID of the auditory cue was  
260 adjusted to match the horizontal eccentricity of the visual target. To ensure that changes  
261 in performance were not due to the spatial localization information provided by the  
262 localized auditory cue, performance between *auditory-static* and *auditory-moving*  
263 conditions was compared (see Figure 4B, note that in both cases sounds were co-  
264 localized with the visual target). Overall, performance accuracy increased by 7.9% in the  
265 presence of a moving, co-localized auditory cue compared to the static co-localized cue.  
266 A two-way ANOVA showed significant main effects of target speed ( $F_{(7, 144)}=42.04$ ,  
267  $p<0.001$ ) and auditory motion ( $F_{(1, 144)}=15.52$ ,  $p<0.001$ ), and no significant interaction  
268 between them ( $F_{(7, 144)}=0.35$ ,  $p>0.9$ ). Thus, in contrast with non-co-localized auditory  
269 cues, where the improvement conferred by congruent motion was small and non-  
270 significant, with spatially co-localized auditory cues there was a significant improvement  
271 in visual performance.

272 We again analyzed reaction times in correct trials and found that response times  
273 decreased from 930 ms in the visual only condition to 827 ms in auditory-static trials, and  
274 further to 775 ms in auditory-moving trials. Neither the difference from visual only to  
275 auditory-static ( $F_{(1, 61)}=3.22$ ,  $p=0.07$ ) nor from auditory-static to auditory-moving ( $F_{(1, 63)}=2.33$ ,  $p=0.13$ ) reached statistical significance in our sample, although the difference  
276 between visual-only and auditory-moving was significant ( $F_{(1, 61)}=20.10$ ,  $p<0.001$ ).  
277 Therefore, the accuracy differences due to auditory motion, observed in the main  
278 analyses, cannot be attributed to speed-accuracy tradeoffs.  
279

280 We replicated the co-localized auditory stimulus experiment with a spectrally  
281 richer auditory stimulus (broadband noise filtered between 200 Hz and 12 kHz), and in  
282 which auditory localization information was conveyed via both interaural level  
283 differences (ILD) and interaural time differences (ITD). Although the localization  
284 information was increased in this condition, resulting in a higher baseline performance,  
285 there was still a remarkable improvement in performance for a congruently moving  
286 auditory cue compared to a static cue (see Supplemental data). Thus, whereas better  
287 auditory stimulus localization may result in a global effect of cross-modal facilitation, our  
288 initial findings indicate that even a relatively coarse auditory motion cue is enough to  
289 provide a significant extra benefit to the detection of moving objects during observer  
290 motion.

### 291 *Auditory localization*

292 Since the baseline and experimental conditions in these experiments both  
293 contained an auditory stimulus presented from the same location, it is unlikely that the  
294 cross-modal benefit reported was due to a spatial cueing effect of the sound. Yet, to  
295 ensure that the increased performance in the auditory-moving trials was not due to  
296 increased localization information provided by the moving auditory cues, we constructed  
297 an auditory localization control test. Three subjects who participated in the main  
298 experiments were presented with an approaching, receding or static auditory cue  
299 (identical to those used in the co-localized auditory cue experiment) localized to one of 9  
300 locations in front of the observer, evenly spaced in  $2.5^\circ$  increments and with no elevation.  
301 After the sound was played, nine vertical bars matching the possible sound source  
302 locations were shown on the screen, and observers were asked to report which one

303 corresponded to the sound origin. We measured the distribution of errors for each sound  
304 condition (see supplemental Figure S2-A). The mean absolute errors were 3.04, 3.06 and  
305 3.04° for receding, static and approaching auditory stimuli, respectively. A Levene test  
306 of variance showed that there were no significant differences in the distributions of errors  
307 for the three cue types for these subjects ( $F_{(2,1733)}=0.06$ ,  $p=0.94$ ). The improved  
308 performance in the auditory-moving condition cannot, therefore, be attributed to an  
309 improved ability to localize the sound source in these trials. A similar pattern arose with  
310 the auditory stimuli containing richer localization cues (Supp. Data, and Fig S2-B).

### 311 *Correlation to visual performance*

312 The strength of multisensory integration has been found to vary as a function of  
313 the accuracy within each modality [45, 46]. We were interested in determining whether  
314 this auditory-based enhancement in visual motion was more effective in observers that  
315 performed poorly on the visual task. To test this, we performed a one-tailed Pearson  
316 correlation test for a negative correlation between the gain due to the moving auditory  
317 cue (auditory-moving relative to auditory-static performance) and performance on the  
318 visual-only task. Note that correlations were made relative to performance on the visual  
319 only condition, so that regression to the mean would not artificially contribute to a  
320 correlation (e.g. a noisy auditory-static data point might cause a noisy measure of  
321 auditory improvement, but this would not be correlated to variations in visual-only  
322 performance). In the non-co-localized condition, where cross-modal benefit was very  
323 modest and not statistically significant, there was no significant relationship between  
324 baseline visual performance and cross-modal gain with moving auditory cues ( $R^2=-0.09$ ,  
325  $p=0.2$ , with  $R^2$  sign assigned based on the  $r$ -value, and indicating a negative correlation).

326 In contrast, with co-localized auditory cues, the correlation was considerably stronger and  
327 statistically significant ( $R^2=-0.37$ ,  $p=0.04$ ). The significant negative correlation shows  
328 that subjects who performed worse on the visual task benefited more from the auditory  
329 cue.

### 330 **Discussion**

331 This study addressed cross-modal enhancement in the detection of visual object  
332 motion during simulated observer motion. Participants were asked to make a visual  
333 discrimination to identify a moving target sphere amidst a dynamic scene simulating an  
334 observer translating forward. We first showed that the pattern of visual search results  
335 was independent of observer speed, indicating that subjects did not resort to performing  
336 the task on the basis of object motion relative to the observer. This result is consistent  
337 with the hypothesis that scene-relative object motion during simulated forward self-  
338 motion is detected by flow-parsing [7-9], in which observer self-motion is estimated and  
339 subtracted from the flow field.

340 Yet, the main finding to emerge from the present study is that the presence of a  
341 moving auditory cue facilitates parsing out relative object motion from optic flow.  
342 Figure 5 summarizes performance across five auditory conditions (visual only, static and  
343 moving non-co-localized auditory cues, and static and moving co-localized auditory  
344 cues). The cross-modal improvement was not due to the mere presence of a sound, given  
345 that accessory static sounds did not result in any advantage, as compared to visual only  
346 displays (static, non-co-localized auditory condition ( $t(26)=0.18$ ,  $p=0.85$ ); or the static,  
347 localized auditory condition ( $t(26)=0.67$ ,  $p=0.5$ ). Additionally, the spatial localization  
348 provided by the auditory cue did not directly improve subject performance: an ANOVA

349 combining the non-co-localized and co-localized conditions showed a significant effect  
350 of auditory motion ( $F_{(1, 288)}=17.8, p<0.001$ ) and a significant interaction between auditory  
351 motion and co-localization ( $F_{(1, 288)}=3.7, p=0.05$ ), but no effect of co-localization alone  
352 ( $F_{(1, 288)}=0.55, p>0.4$ ). Thus, in our task, simply adding a temporally synchronous, static  
353 auditory stimulus did not improve subject performance by either alerting to the stimulus  
354 onset (as in, e.g., [36]), or by directing the observer's attention to the region of the visual  
355 stimulus containing the target object .

356 ----- Figure 5 near here -----

357       However, for moving auditory cues, spatial coincidence between sound and visual  
358 object proved critical, given that congruently moving sounds significantly enhanced  
359 object motion detection only if spatially co-localized. Interestingly, the cross-modal gain  
360 seen for co-localized, moving auditory cues was negatively correlated to individual  
361 performance levels, such that participants who performed worse visually benefited more  
362 from auditory motion. This trend suggests the possibility that auditory cues may be  
363 especially useful to observers with weak visual abilities and thus could be useful in the  
364 rehabilitation of visual deficits. This finding agrees well with previous indications that  
365 visuo-spatial deficits can be ameliorated by using co-localized accessory acoustic cues  
366 [13, 47]. It also supports the idea that the gain of multisensory integration depends on  
367 prior precision levels in unisensory performance [45, 46].

368       Our results therefore suggest that visual-auditory motion integration is more  
369 effective when both cues are presented in spatially commensurate locations within the  
370 stimulus, as has been suggested as a condition for visual-auditory motion binding [14,  
371 18]. Spatially-dependent cross-modal enhancement has been frequently reported in the

372 literature [48, 49], and often linked to the spatial rule of cross-modal integration derived  
373 from single-cell studies in the Superior Colliculus of several animal species (e.g. [50]).

374         Yet, some important exceptions to this rule have been reported recently (e.g.  
375 [51]). Indeed, the strong effect of auditory co-localization in our data is interesting, given  
376 recent reports of cross-modal improvement in visual search tasks that were obtained with  
377 spatially non-informative sounds [36]. This difference between results is however  
378 difficult to interpret at present, given that these previous studies did not include an  
379 auditory co-localized condition to compare with. An interesting speculation, however, is  
380 that in contrast with previous studies of cross-modal enhancement in visual search, the  
381 participants' task in our study was strongly spatial, thus more likely to benefit from  
382 accurate information about spatial relations.

383         A potential mechanism underlying this spatial selectivity is that the co-localized  
384 auditory cue reduced the search space by directing the observer to the approximate  
385 location of the visual target. This could help reduce effective set-size, and thus perceptual  
386 load, allowing audio-visual integration to be more effective. This explanation is indirectly  
387 supported by previous findings indicating that cross-modal integration under high-  
388 perceptual load conditions is mediated by a serial, attentive, process [38, 39, 52], and  
389 therefore should be more effective in conditions where there are fewer possible auditory-  
390 visual associations. Audio-visual coincidence selection can be enabled in a variety of  
391 ways, such as using sparse visual displays (like in many multisensory enhancement  
392 experiments), or by the saliency and temporal informativeness of the accessory acoustic  
393 cue [36]. We hypothesize that the co-localized cues enable efficient audio-visual motion

394 integration since they restrict the search space so that audiovisual motion integration  
395 becomes more effective.

396         The results presented here suggest that parsing object motion from the perceived  
397 optic flow induced by observer self-motion can be enhanced by the presentation of a  
398 spatially co-localized auditory cue of congruent motion. The use of auditory information  
399 in flow-parsing suggests that flow parsing can be seen as a multisensory process, or at  
400 least it is able to operate on multisensory motion representations. A recent magneto-  
401 encephalography (MEG) study of dynamic connectivity among cortical areas involved in  
402 the visual-only and auditory-motion versions of this task [53] found that the middle  
403 prefrontal cortex (MPFC) strongly and selectively modulated the middle temporal area  
404 (MT+) in the visual-only condition, while in the auditory-visual condition, MPFC  
405 provided feedback to the superior temporal polysensory area (STP), to which both the  
406 auditory cortex and MT+ were functionally connected. These results suggested that in  
407 these two tasks the prefrontal cortex allocates attention to the “target” as whole, and that  
408 the target’s representation shifts from MT+ for a moving visual object when no auditory  
409 information was presented, to STP for a moving visual-auditory object. Taken together  
410 with the results we have presented here, we suggest that flow parsing, previously thought  
411 of as a purely visual process, may use multisensory object representations when detecting  
412 a moving object during observer self-motion, demonstrating that the integration of  
413 motion information across sensory modalities contributes to ecological perception that  
414 occurs at early stages of processing.

## 415 Acknowledgements

416 LMV and FJC were supported by NIH grant RO1NS064100, to LMV. SSF was  
417 supported by grants from the *Ministerio de Ciencia e Innovación* (PSI2010-15426 and  
418 CSD2007-00012), by the *Comissionat per a Universitats i Recerca del DIUE* (SRG2009-  
419 092), and by the European Research Council (StG-2010 263145).

420 We thank Gerald Kidd and Chris Mason for their helpful suggestions and for  
421 generously making available to us the resources of the Sound Field Laboratory at Sargent  
422 College, Boston University, supported by grant P30 DC04663. We also thank Franco  
423 Rupcich, Benvy Caldwell and Megan Menard for helping with psychophysical data  
424 collection and subject recruitment, and Leonardo Sassi for developing and implementing  
425 a preliminary version of the object motion task.

426

## 427 References

- 428 [1] Dick M, Ullman S, Sagi D. Parallel and serial processes in motion detection.  
429 Science. 1987;237:400-2.
- 430 [2] McLeod P, Driver J, Crisp J. Visual search for a conjunction of movement and  
431 form is parallel. Nature. 1988;320:154-5.
- 432 [3] Nakayama K, Silverman GH. Serial and parallel processing of visual feature  
433 conjunctions. Nature. 1986 Mar 20-26;320(6059):264-5.
- 434 [4] Duncan J. Target and nontarget grouping in visual search. Perception &  
435 Psychophysics. 1995 Jan;57(1):117-20.
- 436 [5] Gibson JJ. The perception of the visual world. Boston,: Houghton Mifflin 1950.
- 437 [6] Vaina LM. Complex motion perception and its deficits. Current Opinion in  
438 Neurobiology. 1998 Aug;8(4):494-502.
- 439 [7] Rushton SK, Warren PA. Moving observers, relative retinal motion and the  
440 detection of object movement. Curr Biol. 2005 Jul 26;15(14):R542-3.
- 441 [8] Rushton SK, Bradshaw MF, Warren PA. The pop out of scene-relative object  
442 movement against retinal motion due to self-movement. Cognition. 2007  
443 Oct;105(1):237-45.
- 444 [9] Warren PA, Rushton SK. Perception of object trajectory: parsing retinal motion  
445 into self and object movement components. J Vis. 2007;7(11):2 1-11.

- 446 [10] Ross LA, Saint-Amour D, Leavitt VM, Javitt DC, Foxe JJ. Do you see what I am  
447 saying? Exploring visual enhancement of speech comprehension in noisy  
448 environments. *Cereb Cortex*. 2007 May;17(5):1147-53.
- 449 [11] Molholm S, Ritter W, Murray MM, Javitt DC, Schroeder CE, Foxe JJ.  
450 Multisensory auditory-visual interactions during early sensory processing in  
451 humans: a high-density electrical mapping study. *Brain Res Cogn Brain Res*.  
452 2002 Jun;14(1):115-28.
- 453 [12] Vroomen J, de Gelder B. Sound enhances visual perception: cross-modal effects  
454 of auditory organization on vision. *J Exp Psychol Hum Percept Perform*. 2000  
455 Oct;26(5):1583-90.
- 456 [13] Ladavas E. Multisensory-based approach to the recovery of unisensory deficit.  
457 *Ann N Y Acad Sci*. 2008 Mar;1124:98-110.
- 458 [14] Soto-Faraco S, Lyons J, Gazzaniga M, Spence C, Kingstone A. The ventriloquist  
459 in motion: illusory capture of dynamic information across sensory modalities.  
460 *Brain Res Cogn Brain Res*. 2002 Jun;14(1):139-46.
- 461 [15] Soto-Faraco S, Spence C, Kingstone A. Cross-modal dynamic capture:  
462 congruency effects in the perception of motion across sensory modalities. *J Exp*  
463 *Psychol Hum Percept Perform*. 2004 Apr;30(2):330-45.
- 464 [16] Soto-Faraco S, Spence C, Lloyd D, Kingstone A. Moving multisensory research  
465 along: Motion perception across sensory modalities. *Current Directions in*  
466 *Psychological Science*. 2004 Feb;13(1):29-32.
- 467 [17] Soto-Faraco S, Kingstone A, Spence C. Multisensory contributions to the  
468 perception of motion. *Neuropsychologia*. 2003;41(13):1847-62.
- 469 [18] Meyer GF, Wuerger SM, Röhrbein F, Zetzsche C. Low-level integration of  
470 auditory and visual motion signals requires spatial co-localisation. *Experimental*  
471 *brain research Experimentelle Hirnforschung Expérimentation cérébrale*  
472 2005:538-47.
- 473 [19] Alais D, Burr D. No direction-specific bimodal facilitation for audiovisual motion  
474 detection. *Brain Res Cogn Brain Res*. 2004 Apr;19(2):185-94.
- 475 [20] Sanabria D, Lupianez J, Spence C. Auditory motion affects visual motion  
476 perception in a speeded discrimination task. *Exp Brain Res*. 2007 Apr;178(3):415-  
477 21.
- 478 [21] Lopez-Moliner J, Soto-Faraco S. Vision affects how fast we hear sounds move.  
479 *Journal of Vision*. 2007;7(12):1-7.
- 480 [22] Cappe C, Thut G, Romei V, Murray MM. Selective integration of auditory-visual  
481 looming cues by humans. *Neuropsychologia* 2009:1045-52.
- 482 [23] Freeman E, Driver J. Direction of visual apparent motion driven solely by timing  
483 of a static sound. *Curr Biol*. 2008 Aug 26;18(16):1262-6.
- 484 [24] Kitagawa N, Ichihara S. Hearing visual motion in depth. *Nature*. 2002 Mar  
485 14;416(6877):172-4.
- 486 [25] Valjamae A, Soto-Faraco S. Filling-in visual motion with sounds. *Acta Psychol*  
487 (Amst). 2008 Oct;129(2):249-54.
- 488 [26] Alink A, Singer W, Muckli L. Capture of auditory motion by vision is represented  
489 by an activation shift from auditory to visual motion cortex. *Journal of*  
490 *Neuroscience* 2008:2690-7.

- 491 [27] Baumann O, Greenlee MW. Neural correlates of coherent audiovisual motion  
492 perception. *Cereb Cortex* 2007;1433-43.
- 493 [28] Wolfe JM. What Can 1 Million Trials Tell Us About Visual Search? American  
494 Psychological Society. 1998;9(1):33-9.
- 495 [29] Verghese P, Pelli DG. The information capacity of visual attention. *Vision Res.*  
496 1992 May;32(5):983-95.
- 497 [30] Duncan J, Humphreys G. Visual search and stimulus similarity. *Psychological*  
498 *Review.* 1989;96:433-58.
- 499 [31] Bravo M, Blake R. Preattentive vision and perceptual groups. *Perception.*  
500 1990;19:515-22.
- 501 [32] Cavanagh P, Arguin M, Treisman A. Effect of Surface Medium on Visual Search  
502 For Orientation and Size Features. *J Exp Psychol Hum Percept Perform.*  
503 1990;16(3):479-91.
- 504 [33] Royden CS, Wolfe JM, Klempen N. Visual search asymmetries in motion and  
505 optic flow fields. *Perception & Psychophysics.* 2001 Apr;63(3):436-44.
- 506 [34] Rushton SK, Bradshaw MF. Visual search and motion- is it all relative? *Spatial*  
507 *Vision.* 2000;14:85-6.
- 508 [35] Rushton SK, Bradshaw MF, Warren PA. The pop out of scene-relative object  
509 movement against retinal motion due to self-movement. *Cognition.* 2006 Oct 25.
- 510 [36] Van der Burg E, Olivers CN, Bronkhorst AW, Theeuwes J. Pip and pop:  
511 nonspatial auditory signals improve spatial visual search. *J Exp Psychol Hum*  
512 *Percept Perform.* 2008 Oct;34(5):1053-65.
- 513 [37] Iordanescu L, Guzman-Martinez E, Grabowecky M, Suzuki S. Characteristic  
514 sounds facilitate visual search. *Psychon Bull Rev.* 2008 Jun;15(3):548-54.
- 515 [38] Fujisaki W, Koene A, Arnold D, Johnston A, Nishida S. Visual search for a target  
516 changing in synchrony with an auditory signal. *Proc Biol Sci.* 2006 Apr  
517 7;273(1588):865-74.
- 518 [39] Alsius A, Navarra J, Campbell R, Soto-Faraco S. Audiovisual integration of  
519 speech falters under high attention demands. *Curr Biol.* 2005 May 10;15(9):839-  
520 43.
- 521 [40] Brainard DH. The Psychophysics Toolbox. *Spat Vis.* 1997;10(4):433-6.
- 522 [41] Pelli DG. The VideoToolbox software for visual psychophysics: transforming  
523 numbers into movies. *Spat Vis.* 1997;10(4):437-42.
- 524 [42] McLeod P, Driver J, Dienes Z, Crisp J. Filtering by Movement in Visual Search.  
525 *Journal of Experimental Psychology: Human perception and performance.*  
526 1991;17(1):55-64.
- 527 [43] Barker LE, Luman ET, McCauley MM, Chu SY. Assessing equivalence: an  
528 alternative to the use of difference tests for measuring disparities in vaccination  
529 coverage. *Am J Epidemiol.* 2002 Dec 1;156(11):1056-61.
- 530 [44] Barker L, Rolka H, Rolka D, Brown C. Equivalence testing for binomial random  
531 variables: Which test to use? *American Statistician.* 2001 Nov;55(4):279-87.
- 532 [45] Alais D, Burr D. The ventriloquist effect results from near-optimal bimodal  
533 integration. *Curr Biol.* 2004 Feb 3;14(3):257-62.
- 534 [46] Ernst MO, Banks MS. Humans integrate visual and haptic information in a  
535 statistically optimal fashion. *Nature.* 2002 Jan 24;415(6870):429-33.

- 536 [47] Frassinetti F, Bolognini N, Bottari D, Bonora A, Ladavas E. Audiovisual  
537 integration in patients with visual deficit. *J Cogn Neurosci*. 2005 Sep;17(9):1442-  
538 52.
- 539 [48] Frassinetti F, Bolognini N, Ladavas E. Enhancement of visual perception by  
540 crossmodal visuo-auditory interaction. *Experimental Brain Research*  
541 *Experimentelle Hirnforschung Experimentation Cerebrale*. 2002 Dec;147(3):332-  
542 43.
- 543 [49] Bolognini N, Frassinetti F, Serino A, Ladavas E. "Acoustical vision" of below  
544 threshold stimuli: interaction among spatially converging audiovisual inputs. *Exp*  
545 *Brain Res*. 2005 Jan;160(3):273-82.
- 546 [50] Stein B, Meredith M. *The Merging of the Senses*. Cambridge, Massachusetts: The  
547 MIT Press 1993.
- 548 [51] Murray MM, Molholm S, Michel CM, Heslenfeld DJ, Ritter W, Javitt DC, et al.  
549 Grabbing your ear: rapid auditory-somatosensory multisensory interactions in  
550 low-level sensory cortices are not constrained by stimulus alignment. *Cereb*  
551 *Cortex*. 2005 Jul;15(7):963-74.
- 552 [52] Talsma D, Senkowski D, Soto-Faraco S, Woldorff MG. The multifaceted  
553 interplay between attention and multisensory integration. *Trends in cognitive*  
554 *sciences*. Sep;14(9):400-10.
- 555 [53] Vaina L, Calabro F, Lin F, Hamalainen M. Long-Range Coupling of Prefrontal  
556 Cortex and Visual (MT) or Polysensory (STP) Cortical Areas in Motion  
557 Perception. *BIOMAG2010, IFBME Proceedings Series*, Springer Verlag IFBME.  
558 2010.
- 559

560 **Figure captions**

561 **Figure 1** (A) Stimulus display during simulated forward translation (motion vectors  
562 indicated by arrows) with one object moving independently within the scene (indicated  
563 by a bold arrow). (B) Zenithal view of the stimulus layout. As the observer (triangle)  
564 moves forward (3 cm/s), the target (filled circle) moves either forward or back within the  
565 scene (open circles). (C) Amplitude envelope of the auditory cues in the visual only  
566 (top), auditory-static (middle) and auditory-moving (bottom) cases.

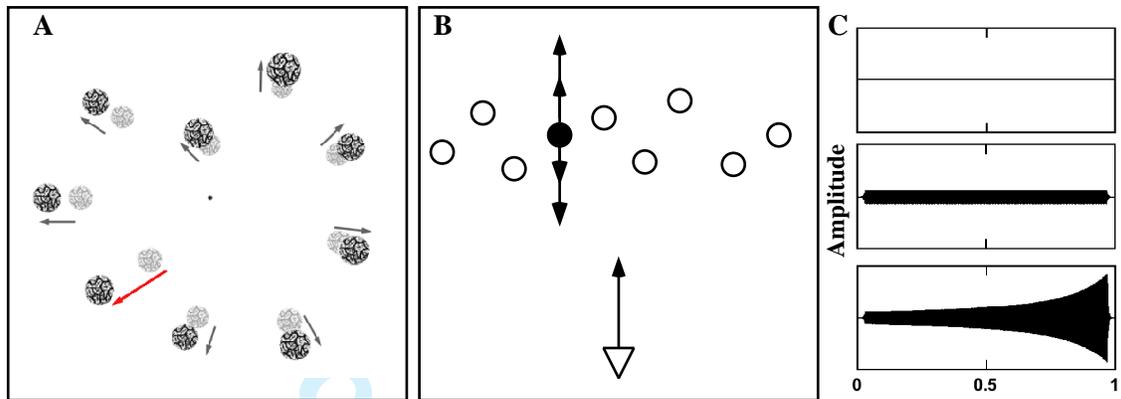
567 **Figure 2** Performance on the visual-only condition for 18 subjects. Error bars are s.e.m  
568 across subjects. Negative speeds refer to receding targets and positive speeds to looming  
569 targets, relative to scene motion. The horizontal line indicates chance performance (25%  
570 correct).

571 **Figure 3** Performance accuracy on the visual-only condition for observer speeds of 3  
572 cm/sec (filled circles) and 5 cm/sec (unfilled circles). Data from the 10 subjects who  
573 participated in both conditions is shown. Arrows indicate the speeds at which an object  
574 would appear stationary on the screen (observer velocity equal to target object velocity).  
575 Error bars are s.e.m. across subjects.

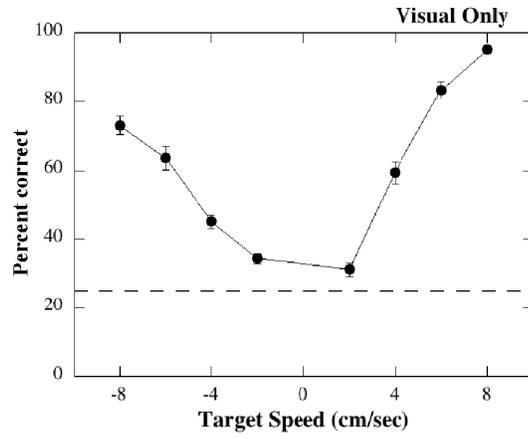
576 **Figure 4** Performance accuracy with (A) a non-co-localized auditory cue and (B) a  
577 spatially co-localized auditory cue, each comparing moving auditory to static auditory  
578 conditions. Error bars are s.e.m. across subjects.

579 **Figure 5** Inter-participant mean motion detection performance across conditions (pooled  
580 across visual target speeds) for the 3 cm/sec observer speed condition. Error bars are  
581 s.e.m. across subjects.

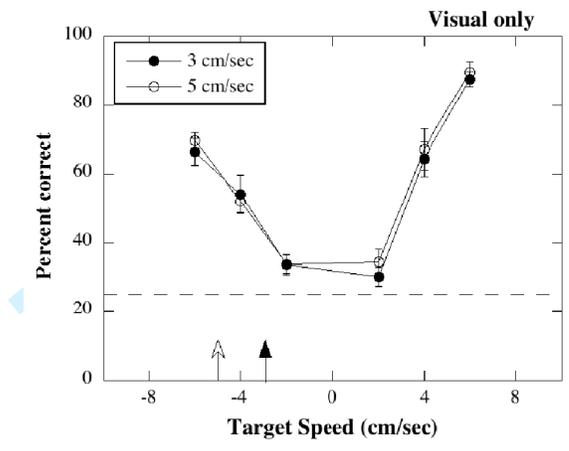
582



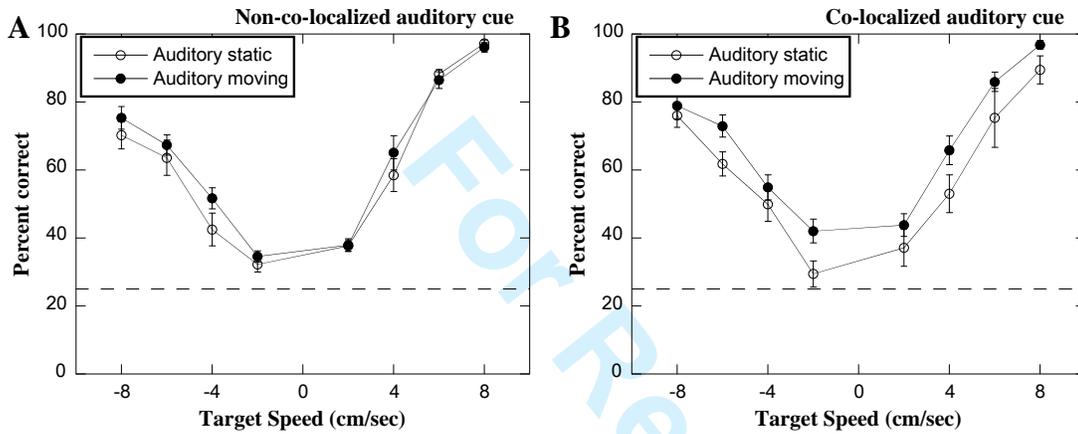
For Review Only

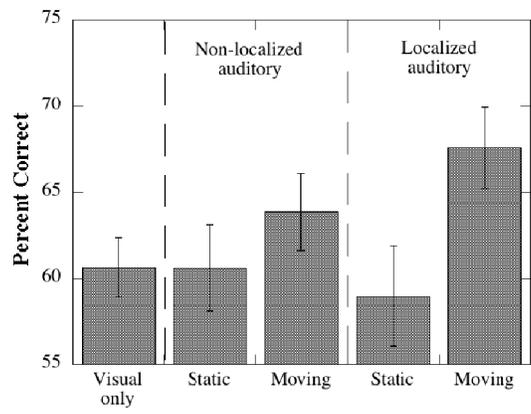


Review Only



Review Only





Review Only