

Effects of auditory reliability and ambiguous visual stimuli on auditory spatial discrimination

Madeline S. Cappelloni^{1,2,3}, Sabyasachi Shivkumar^{3,4}, Ralf M. Haefner^{3,4}, Ross K. Maddox^{1,2,3,5}

¹ *Biomedical Engineering, University of Rochester, Rochester, NY 14627, USA*

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

³ *Center for Visual Science, University of Rochester, Rochester, NY 14627, USA*

⁴ *Brain and Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA*

⁵ *Neuroscience, University of Rochester, Rochester, NY 14627, USA*

ABSTRACT

The brain combines information from multiple sensory modalities to interpret the environment. Multisensory integration is often modeled by ideal Bayesian causal inference, a model proposing that perceptual decisions arise from a statistical weighting of information from each sensory modality based on its reliability and relevance to the observer's task. However, ideal Bayesian causal inference fails to describe human behavior in a simultaneous auditory spatial discrimination task in which spatially aligned visual stimuli improve performance despite providing no information about the correct response. This work tests the hypothesis that humans weight auditory and visual information in this task based on their relative reliabilities, even though the visual stimuli are task-uninformative, carrying no information about the correct response, and should be given zero weight. Listeners perform an auditory spatial discrimination task with relative reliabilities modulated by the stimulus durations. By comparing conditions in which task-uninformative visual stimuli are spatially aligned with auditory stimuli or centrally located (control condition), listeners are shown to have a larger multisensory effect when their auditory thresholds are worse. Even in cases in which visual stimuli are not task-informative, the brain combines sensory information that is scene-relevant, especially when the task is difficult due to unreliable auditory information.

1 I. INTRODUCTION

2 When we navigate our surroundings, we encounter sensory information from multiple
3 modalities. Combining complementary information across sensory modalities often helps us
4 construct a more accurate percept of the world. In contrast, combining conflicting or irrelevant
5 sensory information can lead to perceptual errors. In order to optimize perceptual accuracy, the
6 brain must determine whether to combine information across sensory modalities, and if so, how to
7 weigh each sensory modality. Formally, the notion of reliability weighting is described by Bayesian
8 models of cue combination—forced integration (Ernst and Banks, 2002) and more recently causal
9 inference (Körding et al., 2007). In these models, each cue is treated as a measurement of the
10 stimulus with a Gaussian distribution of the likelihood of the stimulus based on that measurement.
11 The multisensory measurement is then a combination of unisensory measurements weighted by the
12 inverse of their relative variances, such that a narrower likelihood distribution will have more
13 influence on the combined percept. Importantly, the causal inference model adds another layer of
14 inference to this model, in which the degree of cue integration depends on the probability that both
15 measurements actually arose from the same event in the world (Körding et al., 2007).

16 Bayesian models of multisensory integration are typically tested in tasks in which the subject
17 can use information from multiple modalities to determine the correct response; for example, an
18 audiovisual localization task in which the subject is asked where a noise and light occurred (Körding
19 et al., 2007). Under good visual conditions, this task gives rise to the “ventriloquist effect”, a bias of
20 auditory location towards the visual stimulus (Howard and Templeton, 1966). However, when the
21 visual stimulus gets blurrier and harder to localize relative to the auditory stimulus, the apparent
22 visual bias weakens or even manifests as an auditory bias of perceived visual location (Alais and
23 Burr, 2004). This demonstrates that the ventriloquist effect is truly a bias of both visual and auditory
24 stimuli towards each other with the magnitude of the bias determined by the relative reliability of
25 each modality. Importantly, in this and other tasks described by the model, there is only one
26 stimulus in each sensory modality, both of which are informative about the correct response. In this
27 scenario, it is optimal for the brain to use multisensory integration to improve its judgment and
28 behavioral performance.

29 Previously, we extended the classical work by increasing the number of visual and auditory
30 cues and found a multisensory effect of visual stimuli on auditory spatial processing (Cappelloni et
31 al., 2019) even when those visual cue did *not* contain any task-relevant information. We asked
32 listeners to perform a concurrent auditory spatial discrimination task in which random visual stimuli
33 were either spatially aligned with two symmetrically separated auditory stimuli or both collocated in
34 the center of the screen, and found a performance benefit when the visual stimuli were spatially
35 aligned. This audiovisual effect goes against the traditional conception of multisensory integration as
36 a mechanism for the optimal combination of information from the environment (Ernst and
37 Bühlhoff, 2004). The benefit provided by the spatially aligned visual stimuli is also not explained by
38 an ideal Bayesian observer (whose response should be invariant to the locations of the task-
39 uninformative visual stimuli), and is counterintuitive in that the visual stimuli do not help to
40 determine the correct response and must instead benefit the listener through process not part of the
41 ideal Bayesian observer model.

42 Here we test the hypothesis that the brain weighs auditory and visual stimuli by their relative
43 reliabilities even in the case where the visual stimuli do not provide any information about the
44 correct response and would be ignored by an ideal observer. We modulated the reliability of the
45 auditory stimuli by changing their duration, with longer auditory stimuli being more reliable. We
46 found that the benefit provided by the visual stimuli is larger where subjects had poor auditory
47 thresholds. Our results replicate those of our previous study (Cappelloni et al., 2019) and further
48 investigate the ways in which scene-relevant but task-uninformative stimuli can shape perception
49 providing constraints for future theoretical models.

50 **II. METHODS**

51 **A. Participants**

52 Participants (16 female, 4 male; ages ranging between 18 and 31, mean 21.5 +/- 3 years) with
53 normal hearing (thresholds 20 dB HL or better at 500-8000 Hz) and normal vision (self-reported)
54 gave written informed consent. They were compensated for the full duration of time spent in the
55 lab. Research was performed in accordance with protocol approved by the University of Rochester
56 Research Subjects Review Board.

57 **B. Stimuli**

58 Auditory stimuli were pink noise tokens and harmonic tone complexes with matching
59 spectral envelopes bandlimited to 220–4000 Hz. Stimuli were generated and localized by HRTFs
60 from the CIPIC library using interpolation from python's expyfun library as in (Cappelloni et al.,
61 2019), with the notable difference that we generated the pink noise tokens and harmonic tone
62 complexes to be three durations, 100 ms, 300 ms, and 1 s. Auditory stimuli were presented at a
63 24414 Hz sampling frequency and 65 dB SPL level from TDT hardware (Tucker Davis
64 Technologies, Alachua, FL) over ER-2 insert earphones (Etymotic Research, Elk Grove Village, IL).

65 Visual stimuli were regular polygons of per-trial random number of sizes and color. They
66 were inscribed within a 1.5° diameter circle. Colors were chosen to have uniform saturation and
67 luminance, with the two stimuli in each trial having opposite hue as in (Cappelloni et al., 2019).
68 Visual stimuli had the same onset and offset times as the auditory stimuli and thus matched their
69 duration. To prevent overlap they were presented in alternating frames (Blaser et al., 2000) on a
70 monitor with a 144 Hz refresh rate.

71 **C. Task**

72 Each trial began when the subject fixated on a white dot in the center of the screen,
73 confirmed with an eye tracking system (EyeLink 1000, SR Research). Then all four auditory and
74 visual stimuli were presented concurrently for the duration of the trial (100 ms, 300 ms, or 1000 ms).
75 After stimulus presentation, subjects were asked to respond with what side the tone was on by
76 pressing a button. There were two visual conditions: one in which the visual stimuli were spatially
77 aligned with the auditory stimuli and one in which the visual stimuli were collocated in the center of
78 the screen.

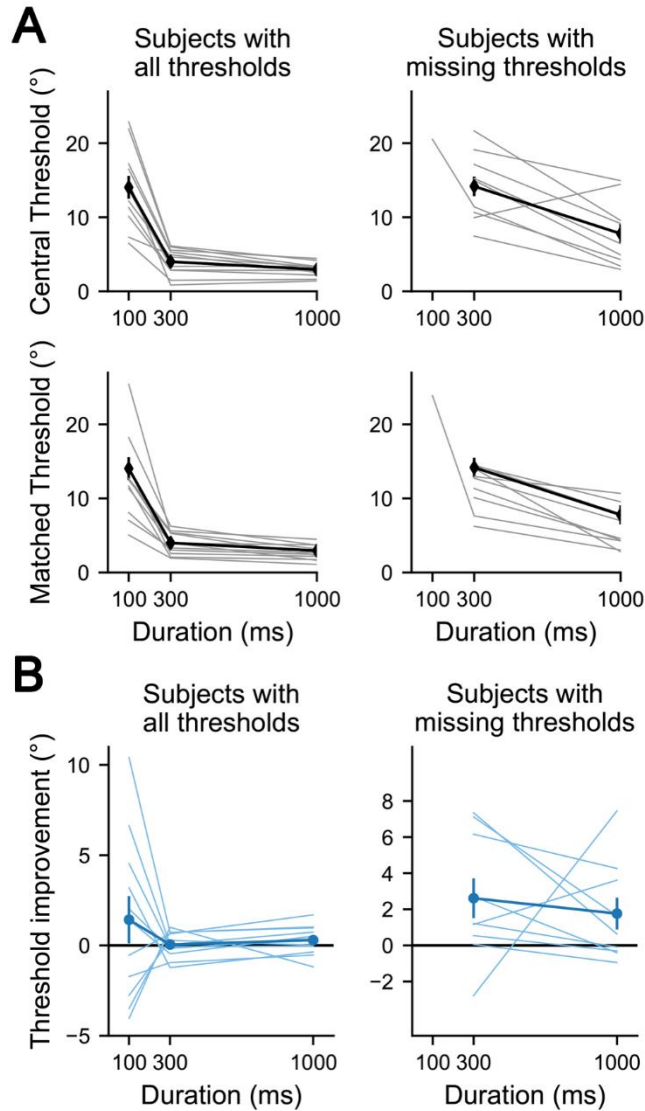
79 We presented trials according to weighted one up one down adaptive tracks that adjusted the
80 separation of the two sounds (Kaernbach, 1991). Separations were adjusted on a log scale such that
81 separation increased by a factor of 2 when the participant responded incorrectly and decreased by a
82 factor of $2^{1/3}$ when they responded correctly. Each track had 130 trials and began at a starting
83 separation of 10° azimuth. For each track, we randomized the number of trials with the tone on the
84 left and right. There were six tracks, three durations by two visual conditions, that were randomly
85 interleaved.

86 **D. Analysis**

87 In order to obtain 75% thresholds we averaged the level at each reversal (skipping the first six
88 reversals). Threshold improvement is defined as the difference between the separation thresholds of
89 the two visual conditions (central – matched). We resampled the reversals with replacement to
90 determine the significance of each threshold improvement (positive or negative respectively – less
91 than 2.5% or greater than 97.5% of resampled threshold improvements less than zero). We
92 performed linear regression on the central threshold vs. threshold improvement data and computed
93 95% confidence intervals using the Python seaborn package (Michael Waskom et al., 2017). We also
94 fit a linear mixed effects model to the data using Python's statsmodels package (Seabold and
95 Perktold, 2010). We fit the thresholds with a random intercept model such that each subject is

- 96 assigned an intercept to control for between subject variance. We considered categorical visual
- 97 condition, duration, and interaction of visual condition and duration as fixed effects in the model.

98 III. RESULTS



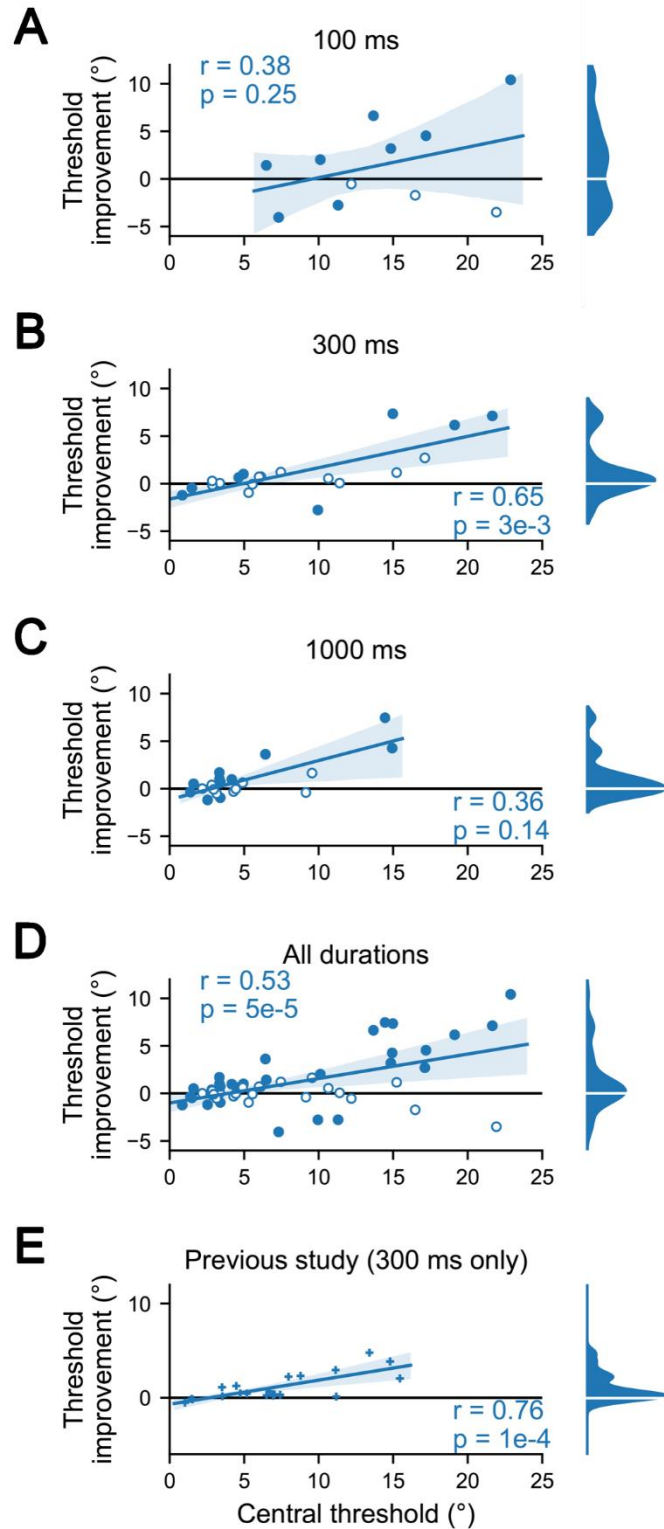
99

100 Fig. 1. (Color Online). A. Thresholds for each duration in the central visual condition (top) and
101 matched visual condition (bottom). Many subjects had missing thresholds (too large to measure
102 accurately) in one or both visual conditions at 100 ms and are plotted in the right column (n=9)
103 while the remainder are plotted on the left (n=11). B. Improvements in threshold for the two groups
104 of subjects: those who could perform the task at all durations (left), those had one or both
105 thresholds missing at 100 ms (right).

106 Subjects improved their task performance, indicated by a decrease in threshold,
107 asymptotically in both visual conditions as the duration of the auditory stimuli increased; however,
108 there was considerable variation in subject performance. Only 11 of 20 subjects were able to
109 perform the auditory discrimination task at the shortest duration such that we could calculate a
110 separation threshold (Figure 1A). Subjects in this group had a large decrease in threshold between
111 100 ms and 300 ms, but did not improve further for 1000 ms stimuli. For the remainder of the
112 subjects who had thresholds too big to calculate in either or both 100 ms conditions, they improved
113 their threshold between 300 ms and 1000 ms. In a linear mixed effect model of all subjects
114 combined, only duration ($p=9.5 \times 10^{-9}$) had a significant effect on threshold. Visual condition and the
115 interaction of visual condition and duration did not have a significant effect on threshold.

116 We defined “threshold improvement” as the difference between the central and matched
117 visual conditions and used it to measure the size of the visual benefit (Figure 1B). Differences in
118 individual auditory spatial processing ability indicate that auditory reliability was not uniform within
119 a given duration condition.

120



122 Fig. 2. (Color Online). Linear regression of threshold improvements against central threshold. Solid
123 markers indicated significant differences between the two visual conditions based on within subject
124 variation ($\alpha = 0.05$, uncorrected). Open markers indicate no significant effect of visual stimuli.
125 Also shown are marginal kernel density estimates (excluding tails beyond which there is less than
126 0.5% of the mass). A-C. Separate regressions for each duration condition. D. Regression for
127 threshold improvements regardless of duration. E. Regression of data from our previous study
128 (Cappelloni et al., 2019).

129 In order to compensate for individual differences we used the separation threshold in the
130 central condition as a measure of auditory reliability for each duration condition (Figure 2). Larger
131 thresholds indicated that the task was more difficult and the individual's auditory reliability was likely
132 worse. Pooling measurements across durations, we found a linear relationship between the threshold
133 improvement and central threshold (Figure 2D, $r=0.53$, $p=5 \times 10^{-5}$). A similar trend is shown when
134 plotting data from our previous study (Cappelloni et al., 2019) (only using 300 ms stimuli) on the
135 same axes (Figure 2E, $r=0.76$, $p=10^{-4}$).

136 IV. DISCUSSION

137 We found that performance in the auditory task correlated with the benefit subjects receive from
138 task-uninformative visual stimuli. Listeners experienced the most benefit when the auditory task was
139 difficult for them (large central threshold). In contrast, individuals who did well in the auditory task
140 (small central threshold) experienced no benefit or even a slight decrement.

141 It should be noted that as the central threshold gets larger, the stimuli become more peripheral
142 where spatial acuity is worse (Hafter and Maio, 1998; Maddox et al., 2014; Middlebrooks and Onsan,
143 2012; Mills, 1958), compounding listeners' worse auditory reliability. A similar decrease may also
144 occur in visual location reliability. Extending the duration improves listener's thresholds in both

145 visual conditions, which can be explained by an improvement in the auditory reliability, suggesting
146 that our duration manipulation does roughly correlate with reliability.

147 This work replicates our previous finding that task uninformative but spatially aligned stimuli
148 benefit auditory spatial discrimination. In the previous study, the visual stimuli preceded auditory
149 stimuli by 100 ms, whereas in this study, their onsets were all concurrent. Additionally, we previously
150 used sigmoidal fits to determine threshold instead of adaptive tracks. This suggests that the visual
151 benefit is robust to small audiovisual asynchronies and changes in probabilistic distribution of
152 stimuli across space (adaptive tracks will lead to more non-uniform priors).

153 Although the visual benefit was larger where subjects showed worse auditory performance, and
154 duration had a significant effect on task difficulty, we did not see a significant effect of changing the
155 duration on visual benefit. This is mainly due to the wide range of auditory spatial processing
156 abilities among the subjects. Because of differences in auditory spatial processing, the effect of
157 duration on visual benefit was inconsistent across subjects. Subjects could be divided roughly into
158 two groups with two different patterns of thresholds, those who could reliably perform the task at
159 100 ms and those who could not. The former group improved their performance significantly when
160 the duration was extended to 300 ms, but did not further improve when the stimuli were 1000 ms,
161 suggesting that they reached ceiling performance at 300 ms. In contrast, the latter group improved
162 significantly when the stimulus duration extended from 300 ms to 1000 ms and were not at ceiling
163 performance at 300 ms. In both this study and our original experiment (Cappelloni et al., 2019),
164 which only included the 300 ms duration condition, we observed a wide range of auditory
165 thresholds. In addition to simple variability across subjects, some thresholds were missing data
166 points because the monitor on which visual stimuli were displayed could not extend far enough to
167 accurately measure large thresholds. These missing data may have allowed us to better fit a model
168 that could show an effect of changing the stimulus duration on visual benefit if such an effect

169 existed. Without considering effects on the scale of individual subjects, for which we did not have
170 enough data, we could not establish a causal relationship between changing stimulus duration and
171 the visual benefit even though correlations suggest one may exist.

172 It is possible that auditory reliability is the underlying factor driving the relationship between
173 auditory performance and visual benefit, even though our data did not show a clear relationship
174 between duration and visual benefit. If auditory reliability does modulate the effect of task-
175 uninformative visual stimuli, following the spirit of Bayesian causal inference, it would further
176 violate the central assumption that the brain will only integrate information that is relevant to the
177 task. This would point to a broader notion of multisensory perception in which stimuli are
178 integrated based on their reliability in representing the sensory scene, rather than the reliability of
179 information they provide regarding a specific task. Further work is needed to describe the
180 boundaries of what information is integrated in a scene. It is important that such work go beyond
181 traditional paradigms to those that can reveal differences of scene-relevant vs task-informative cues

182 **V. CONCLUSION**

183 We show that listeners gain a larger benefit from task-uninformative visual stimuli in an auditory
184 spatial discrimination task when the auditory task is difficult. Our results are consistent with, but do
185 not confirm, the notion that reliability weighting as described in Bayesian models may occur even
186 when visual stimuli do not carry information about the correct decision in the task. We believe two
187 next steps would clarify the findings of this paper. “Small-n” design in which few subjects are
188 recruited to complete many trials would allow us to understand perception on the level of
189 individuals, rather than generalizing across a diverse population (Smith and Little, 2018). Secondly,
190 we call for an exploration of more complex paradigms with multiple multimodal cues caused by
191 potentially multiple events in the world that provide new and stronger tests of existing models.

192 **ACKNOWLEDGMENTS**

193 The authors wish to acknowledge Sara Fiscella for assisting with data collection.

194

195 Research reported in this publication was supported by the National Institute on Deafness and

196 Other Communication Disorders of the National Institutes of Health under award number

197 R00DC014288.

198 REFERENCES

199 Alais, D., and Burr, D. (2004). “The Ventriloquist Effect Results from Near-Optimal Bimodal

200 Integration,” *Curr. Biol.*, **14**, 257–262. doi:10.1016/j.cub.2004.01.029

201 Blaser, E., Pylyshyn, Z. W., and Holcombe, A. O. (2000). “Tracking an object through feature

202 space,” *Nature*, **408**, 196-.

203 Cappelloni, M. S., Shivkumar, S., Haefner, R. M., and Maddox, R. K. (2019). “Task-uninformative

204 visual stimuli improve auditory spatial discrimination in humans but not the ideal observer,”

205 *PLOS ONE*, **14**, e0215417. doi:10.1371/journal.pone.0215417

206 Ernst, M. O., and Banks, M. S. (2002). “Humans integrate visual and haptic information in a

207 statistically optimal fashion,” *Nature*, **415**, 429–433. doi:10.1038/415429a

208 Ernst, M. O., and Bühlhoff, H. H. (2004). “Merging the senses into a robust percept,” *Trends Cogn.*

209 *Sci.*, **8**, 162–169. doi:10.1016/j.tics.2004.02.002

210 Hafter, E. R., and Maio, J. D. (1998). “Difference thresholds for interaural delay,” *J. Acoust. Soc.*

211 *Am.*, **57**, 181. doi:10.1121/1.380412

212 Howard, I. P., and Templeton, W. B. (1966). *Human spatial orientation*, Human spatial orientation,

213 John Wiley & Sons, Oxford, England, 533 pages.

214 Kaernbach, C. (1991). “Simple adaptive testing with the weighted updown method,” *Percept.*

215 *Psychophys.*,

- 216 Körding, K. P., Beierholm, U., Ma, W. J., Quartz, S., Tenenbaum, J. B., and Shams, L. (2007).
217 “Causal Inference in Multisensory Perception,” PLOS ONE, **2**, e943.
218 doi:10.1371/journal.pone.0000943
- 219 Maddox, R. K., Pospisil, D. A., Stecker, G. C., and Lee, A. K. C. (2014). “Directing Eye Gaze
220 Enhances Auditory Spatial Cue Discrimination,” Curr. Biol., **24**, 748–752.
221 doi:10.1016/j.cub.2014.02.021
- 222 Michael Waskom, Olga Botvinnik, Drew O’Kane, Paul Hobson, Saulius Lukauskas, David C
223 Gemperline, Tom Augspurger, et al. (2017). *mwaskom/seaborn: v0.8.1* (September 2017),
224 Zenodo. doi:10.5281/zenodo.883859
- 225 Middlebrooks, J. C., and Onsan, Z. A. (2012). “Stream segregation with high spatial acuity,” J.
226 Acoust. Soc. Am., **132**, 3896–3911. doi:10.1121/1.4764879
- 227 Mills, A. W. (1958). “On the Minimum Audible Angle,” J. Acoust. Soc. Am., **30**, 237–246.
228 doi:10.1121/1.1909553
- 229 Seabold, S., and Perktold, J. (2010). “Statsmodels: Econometric and Statistical Modeling with
230 Python,” Austin, Texas, 92–96. Presented at the Python in Science Conference.
231 doi:10.25080/Majora-92bf1922-011
- 232 Smith, P. L., and Little, D. R. (2018). “Small is beautiful: In defense of the small-N design,” Psychon.
233 Bull. Rev., , doi: 10.3758/s13423-018-1451-8. doi:10.3758/s13423-018-1451-8
234

