1    **Modality-specific and multisensory mechanisms of spatial attention and expectation**

2

3    Authors: Arianna Zuanazzi[1,2], Uta Noppeney[1,3]

4

5    [1]Computational Neuroscience and Cognitive Robotics Centre, University of Birmingham,

6    UK; [2]Department of Psychology, New York University, New York, NY, USA; [3]Donders

7    Institute for Brain, Cognition and Behaviour, Radboud University, Nijmegen, The

8    Netherlands

9

10

11    Corresponding author: Arianna Zuanazzi, Department of Psychology, New York University,

12    New York, NY, USA, email: az1864@nyu.edu

13

14    Number of words: 10917

15

**Abstract**

16

17 In our natural environment, the brain needs to combine signals from multiple sensory

18 modalities into a coherent percept. While spatial attention guides perceptual decisions by

19 prioritizing processing of signals that are task-relevant, spatial expectations encode the

20 probability of signals over space. Previous studies have shown that behavioral effects of

21 spatial attention generalize across sensory modalities. However, because they manipulated

22 spatial attention as signal probability over space, these studies could not dissociate attention

23 and expectation or assess their interaction.

24 In two experiments, we orthogonally manipulated spatial attention (i.e., task-relevance) and

25 expectation (i.e., signal probability) selectively in one sensory modality (i.e., primary

26 modality) (experiment 1: audition, experiment 2: vision) and assessed their effects on primary

27 and secondary sensory modalities in which attention and expectation were held constant.

28 Our results show behavioral effects of spatial attention that are comparable for audition and

29 vision as primary modalities; yet, signal probabilities were learnt more slowly in audition, so

30 that spatial expectations were formed later in audition than vision. Critically, when these

31 differences in learning between audition and vision were accounted for, both spatial attention

32 and expectation affected responses more strongly in the primary modality in which they were

33 manipulated, and generalized to the secondary modality only in an attenuated fashion.

34 Collectively, our results suggest that both spatial attention and expectation rely on modality-

35 specific and multisensory mechanisms.

36

39

40

41 **Introduction**

42

43 Spatial attention is a top-down mechanism that is critical for the selection of task-relevant

44 information. It facilitates perception (e.g., faster reaction times, greater accuracy) of signals

45 presented at the attended location (Carrasco, 2011; van Ede et al., 2012). By contrast, spatial

46 expectation (signal probability) facilitates perception by encoding the statistical structure of

47 the environment (Summerfield & Egner, 2009; Rohenkohl et al., 2014). In everyday life,

48 spatial attention and expectation are closely intertwined. For instance, observers often

49 allocate attentional resources to locations in space where events are more likely to occur

50 (Summerfield & Egner, 2009; Feldman & Friston, 2010 for further discussion within the

51 predictive coding framework).

52 Importantly, in our natural multisensory environment our brain is constantly exposed to

53 auditory and visual signals. This raises the critical question of whether allocation of attention

54 and encoding of signal probability are performed in a modality-specific fashion or

55 interactively across sensory modalities. Previous research has suggested that spatial attention

56 relies on cognitive resources that are partially shared across sensory modalities (Eimer &

57 Driver, 2001; Wahn & König 2015, 2017). For instance, Spence & Driver (1996)

58 manipulated spatial attention by presenting signals with a higher probability in the attended

59 relative to unattended hemifield in one modality only (i.e., primary modality). They showed

60 behavioral facilitation for signals presented at the attended location not only for the primary

61 modality (e.g., audition) but also for the secondary modality (e.g., vision) in which spatial

62 attention was not explicitly manipulated. Likewise, neuroimaging studies showed increased

63 activations for signals presented at the attended location not only in the primary modality in

64 which attention was manipulated but also in the secondary modality (Eimer & Schröger,

65 1998; Eimer, 1999; Macaluso et al., 2002; Santangelo et al., 2009; Zuanazzi & Noppeney,

66  2019). Crucially, in this past work the attentional effects were greater in the primary than in

67  the secondary modality (Spence & Driver, 1996; Mondor & Amirault, 1998). The attenuated

68  generalization across sensory modalities suggests that attentional resources are not

69  supramodal, but partially shared (Driver & Spence, 1998). However, this past research

70  conflated attention and expectation by manipulating spatial attention via probabilistic spatial

71  cues or changes in signal probability (Posner, 1980; Spence & Driver, 1996, 1997; Macaluso

72  et al., 2002; Kincade et al., 2005; Bressler et al., 2008; Santangelo et al., 2009). Notably, the

73  Posner probabilistic cuing paradigm shifts observers' attention via spatial cues that indicate

74  whether a target is, for instance, more likely to appear in the left or right hemifield. Likewise,

75  manipulating not only categorically whether the cue is valid or invalid but also its validity

76  (e.g., 100% vs 60% valid) (Vossel et al., 2006; Doricchi et al., 2010; Macaluso & Doricchi,

77  2013) does not enable the dissociation of spatial attention and expectation. Observers should

78  allocate their attentional resources more to their left hemifield when presented with a cue that

79  indicates with a probability of 1 rather than 0.6 whether the target is likely to be presented in

80  the left hemifield.

81  Thus, the first question of this study is whether spatial attention and/or expectation generalize

82  across sensory modalities to a similar extent, when they are manipulated independently. In

83  the most extreme case, they may be modality-specific (i.e., no generalization) or amodal (i.e.,

84  complete generalization). A recent neuroimaging study, for instance, suggested that spatial

85  attention relies mainly on frontoparietal cortices for both primary and secondary modalities,

86  while spatial expectations are formed in sensory systems selectively for the primary modality

87  (Zuanazzi & Noppeney, 2019). As a consequence, we would expect spatial signal probability

88  to be encoded selectively for the primary modality and to generalize to a secondary modality

89  only to a limited degree.

90  A second unresolved question is whether the generalization across sensory modalities

91  depends on whether attention and expectation are manipulated in the auditory or visual

92  modalities as primary manipulation modality – i.e., on the direction of cross-sensory

93  generalization. Previous studies of multisensory attention have indeed shown asymmetric

94  multisensory generalization, depending on which modality was manipulated as primary

95  modality (Ward at al., 2000; Greene et al., 2001; Molholm et al., 2007). Moreover, one may

96  expect differences in cross-sensory generalization from vision to audition and vice versa

97  because spatial representations and expectations are encoded differently in audition and

98  vision. Visual and auditory systems encode space via different reference frames (i.e., eye-

99  centered vs head-centered) and representational formats. In the visual system, spatial location

100 is directly encoded in the sensory epithelium and later in a place code, i.e., via retinotopic

101 organization of primary and higher order visual cortices (e.g., Sereno et al., 1995; Maier &

102 Groh, 2009). In the auditory system, spatial locations are computed from binaural and

103 monoaural cues in the brain stem and are represented in a hemifield code in primary auditory

104 cortices (e.g., Lauter et al., 1985; Maier & Groh, 2009). Further, in everyday life under

105 normal lighting conditions, vision usually provides more reliable spatial information than

106 audition and therefore often dominates spatial perception (Spence & Driver, 1997; Aller et

107 al., 2015; Odegaard et al., 2015; Rohe & Noppeney, 2015a, 2015b. 2016, 2018; Aller &

108 Noppeney, 2019; Jones et al., 2019; Meijer et al., 2019). As a result, we would expect the

109 generalization of spatial attention and expectation to depend on whether attention and

110 expectation are manipulated primarily in vision or audition.

111 Third, in everyday life spatial expectations are formed when observers implicitly learn the

112 statistical structure of their multisensory environment such as the probability of signals

113 occurring at a particular location. Because spatial information is encoded less reliably in

114 audition than vision, this learning may be faster in vision than in audition. Thereby spatial

115    expectations may also be affected by multisensory processes in perceptual learning (e.g., Kim

116    et al., 2008; Batson et al., 2011). For instance, previous studies suggested that perceptual

117    learning in temporal discrimination tasks generalizes across sensory modalities (Warm et al.,

118    1975; Nagarajan et al., 1998; Meegan et al. 2000; Bratzke et al., 2012; Bueti et al., 2012,

119    2014; but see Lapid et al., 2009). This study will investigate how observers dynamically form

120    spatial expectations by learning signal probability over time (Crist et al., 1997).

121    In two experiments (within participants), we orthogonally manipulated spatial attention and

122    expectation selectively in one sensory modality as primary modality (experiment 1: audition,

123    experiment 2: vision). Crucially, to dissociate spatial attention and expectation we did not use

124    a probabilistic cuing paradigm. Instead, we manipulated observers' spatial attention in the

125    primary (experiment 1: audition, experiment 2: vision) modality by instructing them to attend

126    and respond to, e.g., auditory targets selectively in their left but not right hemifield. In

127    addition, we manipulated the relative frequency of auditory stimuli in the left (e.g., 30%) and

128    right (e.g., 70%) hemifield. Because observers need to respond only to targets in their left

129    hemifield, they should ideally allocate all attentional resources to this task-relevant hemifield

130    irrespective of stimulus frequency. Moreover, observers had to respond to all stimuli in their

131    secondary (e.g., visual) modality irrespective of the hemifield in which they were presented.

132    These visual stimuli were presented equally often in both hemifields. We then assessed the

133    effects of spatial attention and expectation in the primary modality and how they generalize

134    crossmodally to signals in the secondary sensory modality, in which spatial attention and

135    expectation were not explicitly manipulated (experiment 1: vision, experiment 2: audition).

136    If attention, expectation and decision making rely on modality-specific processing streams

137    (i.e., without any multisensory interplay), we would expect as null-hypothesis that the

138    attention and expectation manipulations in the primary modality would not affect response

139    times to stimuli in the secondary modality. Hence, the alternative hypothesis is that both

6

140 spatial attention and expectations rely on mechanisms that are partially shared across sensory

141 modalities and the generalization of spatial attention and expectation depends on whether

142 they are manipulated in audition or vision as primary modalities (Spence & Driver, 1997;

143 Ward at al., 2000; Greene et al., 2001; Molholm et al., 2007; Aller et al., 2015; Odegaard et

144 al., 2015; Rohe & Noppeney, 2015a, 2015b. 2016, 2018; Aller & Noppeney, 2019; Jones et

145 al., 2019; Meijer et al., 2019).

146

147

148 **Materials and Methods**

149

150 **Participants**

151 Twenty-eight healthy subjects (19 females; mean age = 25.57 years; 24 right-handed)

152 participated in the study (experiment 1 and experiment 2, within participants). The sample

153 size was determined based on previous studies that investigated attention/expectation

154 (Doherty et al., 2005; van Ede et al., 2012; Beck et al., 2014; Rohenkohl et al., 2014) and/or

155 multisensory integration (Spence & Driver, 1996, 1997; Eimer et al., 2004; Santangelo et al.,

156 2008; Krumbholz et al., 2009; Mengotti et al., 2018; Zuanazzi & Noppeney, 2018).

157 All participants had normal or corrected to normal vision and reported normal hearing. All

158 participants provided written informed consent and were naïve to the aim of the study. The

159 study was approved by the local ethics committee of the University of Birmingham (Science,

160 Technology, Mathematics and Engineering (STEM) Ethical Review Committee) and the

161 experiment was conducted in accordance with these guidelines and regulations.

162

163

164

165 **Stimuli and Apparatus**

166 Spatial auditory stimuli of 100 ms duration were created by convolving bursts of white noise

167 (with 5 ms onset and offset ramps) with spatially selective head-related transfer functions

168 (HRTFs) based on the KEMAR dummy head of the MIT Media Lab

169 (http://sound.media.mit.edu/resources/KEMAR.html). Visual stimuli ('flashes') were white

170 discs (radius: 0.88° visual angle, luminance: 196 cd/m2) of 100 ms duration presented on a

171 grey background. Both auditory and visual stimuli were presented at ±10° of horizontal visual

172 angle along the azimuth (0° of vertical visual angle). Throughout the entire experiment, a

173 fixation cross was presented in the center of the screen.

174 Prior to the beginning of the study, participants were tested for their ability to discriminate

175 left and right auditory stimuli on a brief series of 20 trials. They indicated their spatial

176 discrimination response (i.e., 'left' vs 'right') via a two-choice key press (group mean

177 accuracy was 99% ± 0.4% [across subjects mean ± SEM]).

178 During the experiment, participants rested their chin on a chinrest with the height held

179 constant across all the participants. Auditory stimuli were presented at approximately 72 dB

180 SPL, via HD 280 PRO headphones (Sennheiser, Germany). Visual stimuli were displayed on

181 a gamma-corrected LCD monitor (2560 x 1600 pixels resolution, 60 Hz refresh rate, 30" Dell

182 UltraSharp U3014, USA), at a viewing distance of approximately 50 cm from the

183 participant's eyes. Stimuli were presented using Psychtoolbox version 3 (Kleiner et al., 2007;

184 www.psychtoolbox.org), running under Matlab R2014a (Mathworks Inc., Natick, MA, USA)

185 on a Windows machine. Participants' responses were recorded via one key of a small keypad

186 (Targus, USA). Throughout the study, participants' eye-movements and fixation were

187 monitored using Tobii Eyex eyetracking system (Tobii EyeX, Tobii, Sweden, ~60 Hz

188 sampling rate).

189

**Study overview: rationale and analysis strategy**

This study included two experiments. Each experiment conforms to a four-factorial design. Because the two experiments were performed within the same participants, the study as a whole could also be treated as a five factorial within-subject experiment. However, because (1) experiment 1 and 2 were completed on different days, (2) experiment 1 was a replication of our previous study and (3) the understanding of results of five factorial designs is rather complex, we will initially analyse each of the two experiments separately.

The separate analyses of experiments 1 and 2 allow us to address our first question, i.e., whether spatial attention and expectations rely on modality-specific or at least partially shared mechanisms. While experiment 1 is intended to replicate our findings reported in our previous research (Zuanazzi & Noppeney, 2018), experiment 2 is intended to extend them and demonstrate that this pattern of results does not depend on whether audition or vision is used as a primary manipulation modality. Moreover, showing the same profile across the two experiments also resolves the ambiguity of our previous research, in which a smaller spatial expectation effect for the secondary (i.e., visual) modality in experiment 1 could potentially be explained by differences in sensory modality rather than attenuated cross-sensory generalization.

In a second step we will directly combine data from experiment 1 and 2 to address our second question, i.e., whether the spatial expectation effect generalizes differently from audition to vision than from vision to audition. To address this question, we need to compare the expectation effects for auditory and visual stimuli in the attended hemifield between the two experiments (i.e., this question cannot be addressed by any of the two experiments alone).

Please also note that the Design and Procedure section mostly overlaps with that of our previous paper (Zuanazzi & Noppeney, 2018) to enable the reader to quickly compare our different studies and obtain a convergent picture across all results.

215 **Design and Procedure**

216 In two experiments, participants were presented with auditory and visual stimuli in their left

217 and right hemifields. To manipulate spatial attention, they were instructed to respond to

218 stimuli in the primary sensory (e.g., auditory) modality selectively in one (i.e., task-relevant)

219 hemifield and ignore stimuli in the task-irrelevant hemifield. Moreover, we manipulated

220 observers' spatial expectations by presenting stimuli in the primary sensory modality with

221 different probabilities in the task-relevant and irrelevant hemifields. In their secondary (e.g.,

222 visual) modality, observers had to respond to all stimuli that were presented equally often in

223 both hemifields (Fig. 1A and 1B). Experiment 1 investigated the effect of auditory spatial

224 attention and expectation on detection of auditory (i.e., primary modality) and visual (i.e.,

225 secondary modality) targets using a 2 (auditory spatial attention: left vs right hemifield) x 2

226 (auditory spatial expectation: left vs right hemifield) x 2 (stimulus modality: auditory vs

227 visual) x 2 (stimulus location: left vs right hemifield) factorial design. Hence, experiment 1

228 manipulated spatial attention and expectation selectively in audition and assessed their direct

229 effects on auditory stimulus processing and indirect generalization to visual stimuli. In

230 experiment 2 primary and secondary modality were reversed (i.e., primary modality: vision;

231 secondary modality: audition); design and procedural details were otherwise comparable to

232 experiment 1. For the data analysis we pooled over stimulus locations (left/right) leading to a

233 2 (auditory spatial attention: attended vs unattended) x 2 (auditory spatial expectation:

234 expected vs unexpected) x 2 (stimulus modality: auditory vs visual) factorial design.

235 Spatial attention was manipulated for the primary modality as task-relevance, i.e., the

236 requirement to respond to an auditory (experiment 1) or a visual (experiment 2) target in the

237 left vs right hemifield. Prior to each run a cue (duration: 2000 ms) informed the observer

238 whether to respond to targets in either their left or right hemifield.

10

239    Spatial expectation was manipulated as spatial signal probability for signals in the primary

240    modality across experimental sessions that were performed on different days. Auditory (i.e.,

241    primary modality in experiment 1) or visual (i.e., primary modality in experiment 2) signals

242    were presented with a ratio of 2.33/1 (i.e., 70%/30%) in the expected/unexpected hemifield.

243    Observers were not informed about those probabilities but learnt them implicitly.

244    Importantly, spatial attention and expectation were not directly manipulated in the secondary

245    modality, allowing us to assess their cross-sensory generalization. As a result, participants

246    needed to respond to all visual targets that were presented with equal probabilities in their

247    spatial hemifields in experiment 1 (i.e., ratio 1/1 in the expected/unexpected hemifields) (Fig.

248    1A and 1C). Likewise, they had to respond to all auditory targets that were presented with

249    equal probabilities in experiment 2.

250    Each experiment included two sessions (i.e., spatial expectation left vs right on different

251    days). Hence, subjects participated in the two experiments on four days separated by at least

252    2 to a maximum of 10 days: 2 sessions for experiment 1 and 2 sessions for experiment 2 = 4

253    sessions in total for each participant. Each session included 12 attention runs. Runs were of

254    two types: in run type A (Fig. 1A, 1C and 1D) spatial attention and expectation were

255    congruent (i.e., spatial attention was directed to the hemifield with higher stimulus

256    frequency); in run type B spatial attention and expectation were incongruent (i.e., spatial

257    attention was directed to the hemifield with less frequent stimuli). The overall probability to

258    respond (i.e., response probability) was greater when attention and expectation were

259    congruent and directed to the same hemifield (85%, runs of type A) than when they were

260    directed to different hemifields (65%, runs of type B) (Fig. 1D).

261    The order of experiments 1 vs 2 and of expectation sessions (i.e., left vs right) was

262    counterbalanced across participants; the order of attention runs (i.e., left vs right) was

263    counterbalanced within and across participants and the order of stimulus locations (i.e., left vs

264 right) and stimulus modalities (sound vs flash) was pseudo-randomized within each

265 participant. Brief breaks were included after every two attention runs to provide feedback to

266 participants about their performance accuracy (averaged across all conditions) in the target

267 detection task and about their eye-movements (i.e., fixation maintenance).

268 Overall, each experiment included 80 trials x 12 attention runs (6 runs of type A and 6 runs of

269 type B, duration: 3 mins/run) x 2 expectation sessions = 1920 trials in total (and 3840 for the

270 whole study). Specifically, each run type included i. 336 targets presented in the expected

271 hemifield (pooled over left and right) and 144 targets in the unexpected hemifield (pooled

272 over left and right) for the primary modality and ii. 240 targets presented in the expected

273 hemifield and 240 targets in the unexpected hemifield (pooled over left and right) for the

274 secondary modality. For further details see Fig. 1C which shows the absolute number of trials
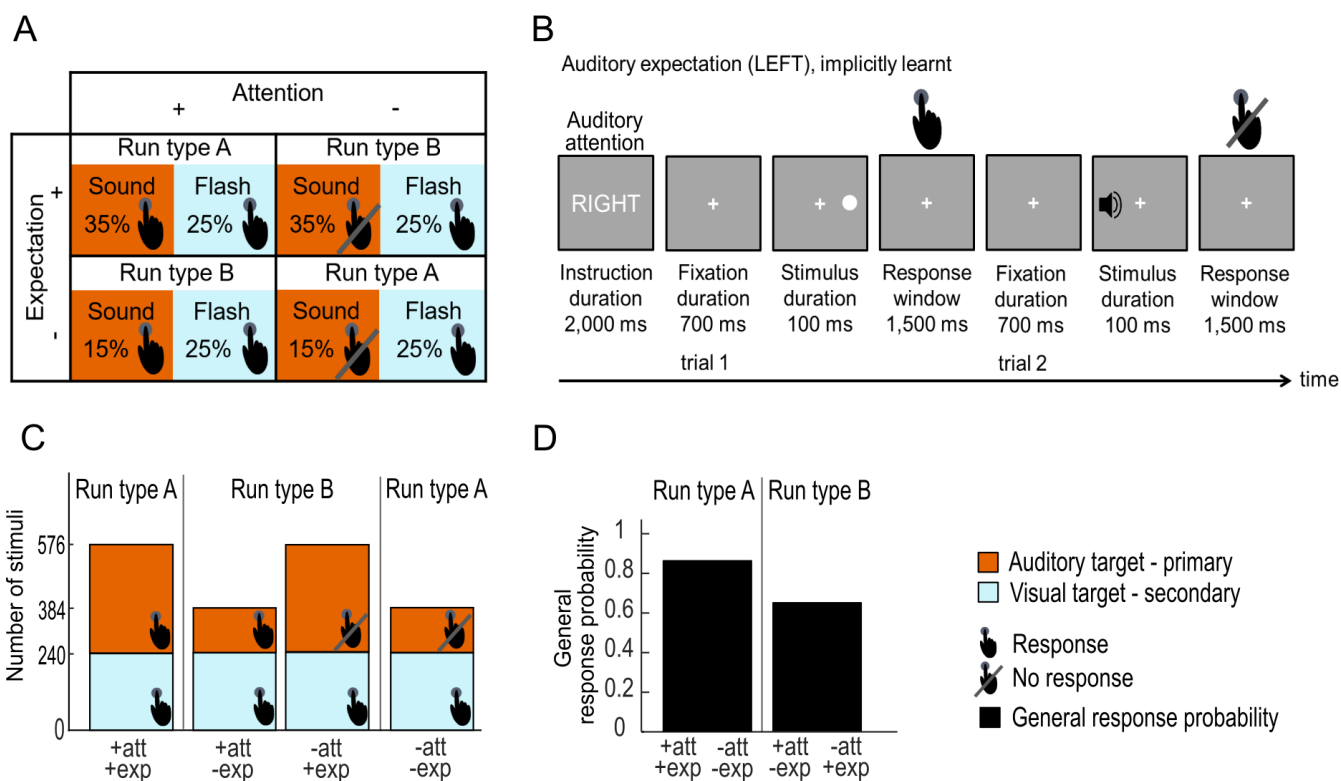
275 for each condition and run type.

276 Each trial (SOA: 2300 ms) included three time windows: i. fixation cross alone (700 ms

277 duration), ii. brief sound or flash (stimulus duration: 100 ms) and iii. fixation cross alone

278 (1500 ms as response interval, see Fig. 1B). Participants responded to the targets in the

279 primary modality presented in the attended hemifield and to all targets in the secondary

280 modality irrespective of hemifield via key press, with their index finger (i.e., the same

281 response for all auditory and visual stimuli) as fast and accurately as possible (Fig. 1B).

282 Prior to each session, participants were familiarized with the stimuli in brief practice runs

283 (with equal spatial signal probability) and trained on target detection performance and

284 fixation (i.e., a warning signal was shown when the disparity between the central fixation

285 cross and the eye-data samples exceeded 2.5 degrees).

286 After the final session (i.e., experiment 1 for 14 participants and experiment 2 for the other 14

287 participants), participants indicated in a questionnaire whether they thought the stimuli in the

288 primary modality were presented more frequently in one of the two spatial hemifields. Ten

289    out of 14 participants in experiment 1 and 11 out of 14 participants in experiment 2 correctly

290    identified the expectation manipulation in the primary modality. Moreover, 13 out of 14

291    participants in experiment 1 and 13 out of 14 participants in experiment 2 correctly reported

292    that stimuli in the secondary modality were presented with equal probabilities across the two

293    hemifields. These data suggest that the majority of participants were aware of the

294    manipulation of signal probability at least at the end of the fourth session.

## Experiment 1 (audition to vision)



**Figure 1**: Design and example trials of experiment 1 (audition to vision)

**A.** Experiment 1: auditory spatial attention and expectation (i.e., signal probability) were manipulated in a 2 (auditory modality – dark orange, vs visual modality – light blue) x 2 (attended hemifield vs unattended hemifield) x 2 (expected hemifield vs unexpected hemifield) factorial design. For illustration purposes, stimulus locations (left/right) were collapsed. Presence vs absence of response requirement is indicated by the hand symbol, spatial signal probability manipulation is indicated by the %. **B.** Experiment 1: example of two trials in a session where auditory stimuli were presented with a probability of 0.7 in the left hemifield and 0.3 in the right hemifield. At the beginning of each run (i.e., 80 trials), a cue informed participants whether to attend and respond to auditory signals selectively in their left or right hemifield throughout the entire run. On each trial participants were presented with an auditory or visual stimulus (100 ms duration) either in their left or right hemifield. They were instructed to respond to auditory stimuli only in the attended hemifield

309 and to all visual stimuli irrespective of the hemifield as fast and accurately as possible with

310 their index finger. The response window was limited to 1500 ms. Participants were not

311 explicitly informed that auditory signals were more likely to appear in one of the two

312 hemifields. Instead, spatial expectation was implicitly learnt within a session (i.e., day).

313 **C.** Experiment 1: number of auditory (dark orange) and visual (light blue) trials in the 2

314 (attended vs unattended hemifield) x 2 (expected vs unexpected hemifield) design (pooling

315 over left/right stimulus location). Presence vs absence of response requirement is indicated by

316 the hand symbol. The fraction of the area indicated by the 'Response' hand symbol pooled

317 over the two bars of one particular run type (e.g., run type A) represents the response related

318 expectation (i.e., general response probability: the overall probability that a response is

319 required on a particular trial); general response probability is greater for run type A (85%),

320 where attention and expectation are congruent, than for run type B (65%), where attention

321 and expectation are incongruent, as indicated in **D**.

322 Note. Design and procedure of experiment 2 were comparable to that of experiment 1, with

323 the only difference that vision was the primary modality and audition was the secondary

324 modality. In other words, in experiment 2 attention and expectation were manipulated

325 selectively in vision.

326

**Data analysis**

**Eye movement: exclusion criteria**

We excluded trials where participants did not successfully fixate the central cross based on a dispersion criterion (i.e., distance of fixation from subject's center of fixation, as defined in calibration trials, > 1.3 degrees for three subsequent samples; Blignaut, 2009). Our eyetracking data confirmed that participants successfully maintained fixation in both experiments with only a small number of trials to be excluded (experiment 1: excluded auditory response trials 1.8% ± 0.5% [across subjects mean ± SEM]; excluded visual response trials 1.7% ± 0.5% [across subjects mean ± SEM]; experiment 2: excluded visual response trials 2.7% ± 1% [across subjects mean ± SEM]; excluded auditory response trials 2.7% ± 0.9% [across subjects mean ± SEM]).

**Response time analysis - separately for experiments 1 and 2**

We initially analysed response times separately for primary and secondary modalities and independently for experiments 1 and 2.

The response time (RT) analysis was limited to trials with RT within the 1500 ms response window and was performed after pooling over stimulus location (left/right).

For the primary modality (i.e., experiment 1 = audition, experiment 2 = vision), subject-specific median RT were entered into a two-sided paired-sample t-test with spatial expectation (expected vs unexpected stimulus) as factor (observers did not respond to targets in the primary modality in the 'unattended' hemifield). Moreover, subject-specific False Alarm rates (FA) for the 'unattended' hemifield were entered into a non-parametric two-sided Wilcoxon Signed Rank tests with spatial expectation (expected vs unexpected stimulus) as factor. We used non-parametric tests, because False Alarms rates are bounded between 0 and 1 and therefore not normally distributed. For the secondary modality (i.e., experiment 1 =

352 vision, experiment 2 = audition), subject-specific median response time were entered into a 2

353 (spatial attention: attended vs unattended stimulus) x 2 (spatial expectation: expected vs

354 unexpected stimulus) repeated measures analysis of variance (ANOVA).

355 For both experiments 1 and 2, the mean hit rates were very high (> 99% in all conditions,

356 Table 1), indicating that participants accurately performed the detection task. Because of the

357 absence of a substantial number of misses, hit rates were not further analyzed.

358

**359 Response time analysis - combined for experiments 1 and 2**

360 To compare effects of primary vs secondary modality, we compared the response times in the

361 attended hemifield (averaged across expected and unexpected hemifields) for auditory and

362 visual stimuli across the two experiments in a 2 (stimulus modality: audition vs vision) x 2

363 (manipulation: primary/direct vs secondary/indirect modality) repeated measures ANOVA.

364 Next, we investigated whether the effect of spatial expectation (i.e., expected vs unexpected)

365 in the attended hemifield (no response was required for stimuli in the primary modalities

366 presented in the unattended hemifield) depended on i. whether targets were presented in the

367 primary or secondary modalities (i.e., the extent to which spatial expectations generalize

368 across the senses) and ii. the multisensory generalization direction, from audition to vision

369 and from vision to audition (i.e., whether spatial expectations generalize differently

370 depending on whether audition or vision is the primary modality). Hence, we first computed

371 the difference in median RT ($\Delta RT_{Exp}$) between unexpected and expected stimuli presented in

372 the attended hemifield (which corresponds to $\Delta RT_{Exp}$ between attended stimuli in run type B

373 and run type A) for targets in each experiment, yielding four conditions: i. auditory targets as

374 primary modality (experiment 1), ii. visual targets as secondary modality (experiment 1), iii.

375 visual targets as primary modality (experiment 2), iv. auditory targets as secondary modality

376 (experiment 2). $\Delta RT_{Exp}$ were entered into a 2 (multisensory generalization direction: audition

17

377    to vision vs vision to audition) x 2 (manipulation: primary/direct vs secondary/indirect

378    modality) repeated measures ANOVA. Please note that the interaction then reflects the

379    difference between targets in the auditory and visual modality.

380

381    **Time course of response times - combined for experiment 1 and 2**

382    Finally, we assessed how these effects of multisensory generalization direction and

383    direct/indirect manipulation evolved over time. For this, we computed the difference in

384    median RT ($\Delta RT_{Exp}$) between unexpected and expected stimuli, as in the previous analysis,

385    but now separately for the first and second half of the experiment (i.e., one half = 430 trials).

386    Each half contained the data from 6 subsequent attention runs (3 runs of type A and 3 runs of

387    type B) for each expectation condition. $\Delta RT_{Exp}$ (or each half) were entered in a 2

388    (multisensory generalization direction: audition to vision vs vision to audition) x 2

389    (manipulation: primary/direct vs secondary/indirect modality) x 2 (time: first vs second half

390    of the experiment) repeated measures ANOVA. Figure 2D shows the across subjects' mean

391    (±SEM) RT separately for each of the 2 attention runs (1 of type A and 1 of type B) (orange

392    and blue circles) as a more fine-grained temporal characterization of the effects of

393    expectation over time. An additional analysis using this more fine-grained temporal division

394    replicated the results reported in this manuscript where we separated the data into halves.

395

396    For all analyses we assessed the assumptions of normality using the Shapiro–Wilk test

397    (Shapiro and Wilk, 1965). When normality was violated, we evaluated the main effects of

398    attention, expectation and their interactions in the factorial design using permutation testing

399    with $2^{28}$ permutations (Nichols & Holmes, 2002). Because in these cases permutation tests

400    replicated the results of the initial ANOVAs, we only report the results of the ANOVAs for

401    consistency.

402 **Results**

403

404 **Generalization of attention and expectation effects across modalities – separately for**

405 **experiment 1 and 2**

406 In experiment 1, participants responded to auditory targets presented in their attended

407 hemifield and to all visual targets. In experiment 2, participants responded to visual targets

408 presented in their attended hemifield and to all auditory targets.

409 We observed qualitatively similar effects across experiments 1 and 2. Table 1 shows RT

410 (across participants' mean ± SEM) for targets in the auditory and visual modalities for the

411 two experiments.

412 For the primary modality, the two-sided paired-sample t-tests showed significantly faster RT

413 in the attended hemifield, when this hemifield was expected than unexpected (experiment 1,

414 auditory modality: $t(27) = -2.83$, $p = 0.009$, Cohen's $d_{av}$ [95% CI] = -0.16 [-0.27, -0.04];

415 experiment 2, visual modality: $t(27) = -9.62$, $p < 0.001$, Cohen's $d_{av}$ [95% CI] = -0.35 [-0.47,

416 -0.23]) (Table 1, Fig. 2A and 2B). Moreover, the Wilcoxon Signed Rank tests showed

417 significantly greater FA in the unattended hemifield, when stimuli in this hemifield were

418 unexpected than expected (experiment 1, auditory modality: W = 52, $p = 0.001$, r [95% CI] =

419 -0.72 [-0.87, -0.45]; experiment 2, visual modality: W = 2, $p < 0.001$, r [95% CI] = -0.99 [-1,

420 -0.98]) (Table 1).

421 For the secondary modality, the 2 (attended vs unattended) x 2 (expected vs unexpected)

422 repeated measures ANOVAs revealed a significant main effect of attention (experiment 1,

423 visual modality: $F(1, 27) = 72.08$, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.73 [0.55, 0.80]; experiment 2,

424 auditory modality: $F(1, 27) = 36.91$, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.58 [0.34, 0.70]). Results

425 showed that participants responded faster to targets presented in their attended than

426 unattended hemifields.

427　Moreover, a significant crossover interaction between attention and expectation was observed

428　(experiment 1, visual modality: $F(1, 27) = 10.09$, $p = 0.004$, $\eta_p^2$ [90% CI] = 0.27 [0.06, 0.46];

429　experiment 2, auditory modality: $F(1, 27) = 44.10$, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.62 [0.40,

430　0.73]) (Table 1, Fig. 2A and 2B). The simple main effects showed that participants responded

431　significantly faster to targets in their attended hemifield, when this hemifield was expected

432　than unexpected (experiment 1, visual modality: $t(27) = -2.81$, $p = 0.009$, Cohen's $d_{av}$ [95%

433　CI] = -0.08 [-0.15, -0.02]; experiment 2, auditory modality: $t(27) = -5.96$, $p < 0.001$, Cohen's

434　$d_{av}$ [95% CI] = -0.14 [-0.19, -0.08]). This significant simple main effect demonstrates that the

435　effects of spatial attention and expectation generalized from primary to secondary modalities,

436　where neither attention nor expectation were explicitly manipulated. By contrast, participants

437　responded significantly more slowly to targets in the secondary modality in the unattended

438　hemifield, when this hemifield was expected than unexpected (experiment 1, visual modality:

439　$t(27) = 2.56$, $p = 0.016$, Cohen's $d_{av}$ [95% CI] = 0.09 [0.02, 0.17]; experiment 2, auditory

440　modality: $t(27) = 5.53$, $p < 0.001$, Cohen's $d_{av}$ [95% CI] = 0.18 [0.10, 0.26]) (Table 1, Fig. 2A

441　and 2B). We suggest that simple main effects for expectation show opposite directions

442　because of response inhibition. In the attended hemifield observers need to respond to the

443　stimuli in the primary modality. Hence, if stimuli from the primary modality are frequent

444　(i.e., expected) in the attended hemifield, observers need to respond on a large percentage of

445　trials. By contrast, in the unattended hemifield observers should not respond to the stimuli in

446　the primary modality. Hence, if stimuli in the primary modality are frequent (i.e., expected)

447　in the unattended hemifield, observers need to inhibit their response on a large percentage of

448　trials. This explanation is also supported by the increase in FA for the primary modality in the

449　unattended hemifield when stimuli in this hemifield are unexpected relative to expected (see

450　above and Fig. 1D). Collectively, the response times and FA rates suggest that, in runs in

451　which observers need to respond to many stimuli, because the stimulus frequency is high in

452      the task-relevant/attended hemifield, observers will make more false alarms to stimuli of the

453      primary modality and respond faster to stimuli of the secondary modality in the unattended

454      hemifield. We can explain this profile in decision making models in which observers need to

455      accumulate evidence to a threshold. An increase in the percentage of trials that require a

456      response may then be reflected either in a shift of the starting point closer to the decisional

457      boundary or in a lower decisional boundary (Gold & Shadlen, 2001).

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477 **Table 1**: Group mean hit rates, false alarm (FA) rate and reaction times (RT) and for each

478 stimulus modality in each condition where a response was given for experiment 1 (primary

479 modality: audition, secondary modality: vision) and experiment 2 (primary modality: vision,

480 secondary modality: audition). In experiment 1, participants responded only to attended

481 auditory targets (and to all visual targets); in experiment 2 participants responded only to

482 attended visual targets (and to all auditory targets). Standard errors (SEM) are given in

483 parentheses.

484

| | Auditory modality | | | | Visual modality | | | |
|---|---|---|---|---|---|---|---|---|
| **Experiment 1** | +att +exp | +att -exp | -att +exp | -att -exp | +att +exp | +att -exp | -att +exp | -att -exp |
| Hit rate (%) (SEM) | 99.5 (0.19) | 99.5 (0.15) | / | / | 99.7 (0.07) | 99.7 (0.07) | 99.3 (0.16) | 99.4 (0.13) |
| FA rate (%) (SEM) | / | / | 4.4 (0.8) | 7.1 (0.1) | / | / | / | / |
| RT (ms) (SEM) | 599.7 (20.1) | 616.2 (19.2) | / | / | 514.4 (16.9) | 522.1 (16.9) | 583.8 (18.5) | 574.4 (18.9) |
| **Experiment 2** | +att +exp | +att -exp | -att +exp | -att -exp | +att +exp | +att -exp | -att +exp | -att -exp |
| Hit rate (%) (SEM) | 99.8 (0.08) | 99.9 (0.03) | 99.2 (0.22) | 99.4 (0.18) | 99.7 (0.09) | 99.7 (0.09) | / | / |
| FA rate (%) (SEM) | / | / | / | / | / | / | 3.2 (0.5) | 8 (1) |
| RT (ms) (SEM) | 527.4 (20.2) | 542.1 (20) | 605.2 (25.7) | 581.3 (24.2) | 526.3 (17.9) | 561.8 (19.9) | / | / |

485

486

487

**The effects of primary vs secondary modality – combined for experiment 1 and 2**

To assess the effect of primary vs secondary modality unconfounded by differences between auditory vs visual modality, we directly compared the response times in the attended hemifield (averaged across expected and unexpected hemifields) for auditory and visual stimuli across the two experiments. The 2 (stimulus modality: audition vs vision) x 2 (manipulation: primary/direct vs secondary/indirect modality) repeated measures ANOVA revealed a significant main effect of stimulus modality ($F(1, 27) = 40.14$, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.60 [0.37, 0.72]), showing faster RT for visual than auditory targets; of manipulation ($F(1, 27) = 151.80$, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.85 [0.74, 0.89]), showing overall faster RT for secondary than primary modality. Moreover, we observed a significant interaction between stimulus modality and manipulation ($F(1, 27) = 6.79$, $p = 0.015$, $\eta_p^2$ [90% CI] = 0.20 [0.02, 0.39]), showing that observers were significantly faster responding to stimuli in the secondary than primary sensory modality predominantly in experiment 1 (primary = auditory, secondary = visual, $t(27) = 11.67$, $p < 0.001$, Cohen's $d_{av}$ [95% CI] = 0.93 [0.63, 1.22]).

Collectively, these results suggest that observers responded faster to stimuli in their secondary modality than in their primary modality. As expected, this effect was more sensitively revealed for auditory stimuli that were associated with slower response times. These effects can be explained by the fact that observers needed to respond to all stimuli in the secondary modality irrespective of hemifield. By contrast, they first needed to discriminate whether signals were presented in the left or right hemifield when responding to stimuli in the primary sensory modality. Because the spatial reliability is lower for auditory than visual signals in our study, these effects were more prominent for auditory stimuli.

**The effects of spatial expectation – combined for experiment 1 and 2**

In the previous analysis we showed that the effects of expectation on response times generalized from the primary to the secondary modality. Next, we directly compared the effects of spatial expectation across the two experiments, in which either audition was the primary and vision the secondary modality or vice versa. The 2 (multisensory generalization direction: audition to vision vs vision to audition) x 2 (manipulation: primary/direct vs secondary/indirect modality) repeated measures ANOVA revealed a significant main effect of manipulation ($F(1, 27) = 18.03$, $p < .001$, $\eta_p^2$ [90% CI] =0.40 [0.16, 0.56]), showing overall greater expectation effects for primary than secondary modalities (dark vs light bars in Fig. 2C). Moreover, we observed a significant main effect of crossmodal generalization direction ($F(1, 27) = 20.71$, $p < .001$, $\eta_p^2$ [90% CI] = 0 .43, [0.19, 0.59]), with a greater expectation effect (i.e., greater $\Delta RT_{Exp}$) for vision to audition than for audition to vision (experiment 2 vs experiment 1 in Fig. 2C). Critically, this generalization effect may be greater from vision to audition than vice versa because observers learn signal probabilities and hence form spatial expectations faster when vision is the primary modality (dark blue bar in Fig. 2C). Alternatively, the expectation effect generalizes more effectively from vision to audition than vice versa (light orange bar in Fig. 2C).

**Time course of the effects of spatial expectation - combined for experiment 1 and 2**

To disentangle between these two possibilities, we investigated how signal probability is learnt over time when audition (experiment 1) or vision (experiment 2) are the primary modality by repeating the previous analysis with the additional factor of time (i.e., first vs second half of experiment). The 2 (multisensory generalization direction: audition to vision vs vision to audition) x 2 (manipulation: primary/direct vs secondary/indirect modality) x 2 (time: first vs second half of the experiment) repeated measures ANOVA performed on

538    $\Delta RT_{Exp}$ revealed significant main effects of multisensory generalization direction ($F(1, 27)$ =

539    22.70, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.46 [0.21, 0.61]), manipulation ($F(1, 27)$ = 13.77, $p <$

540    0.001, $\eta_p^2$ [90% CI] = 0.34 [0.10, 0.51]) as well as a significant interaction between

541    multisensory generalization direction x manipulation x time ($F(1, 27)$ = 11.38, $p$ = 0.002, $\eta_p^2$

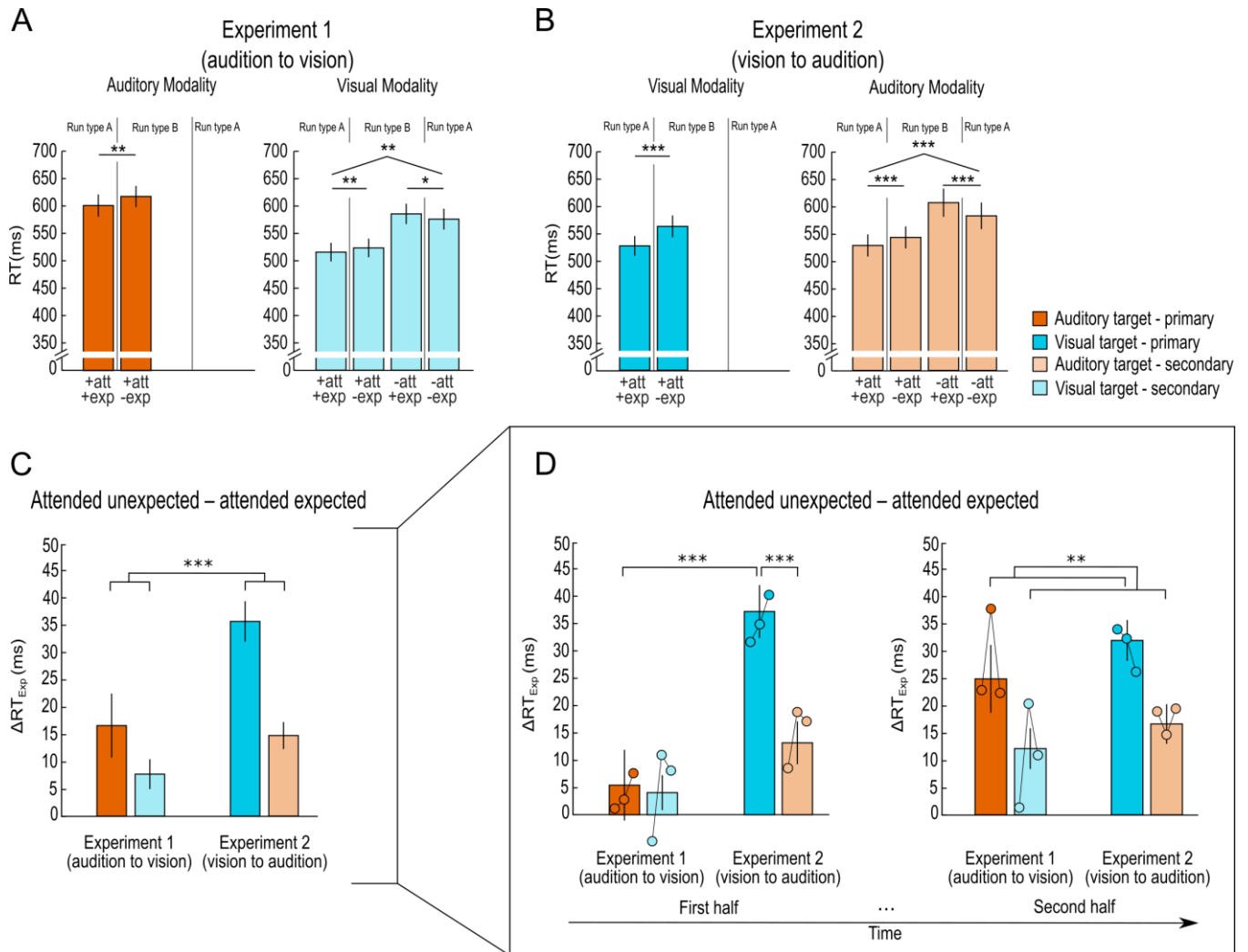542    [90% CI] = 0.29 [0.07, 0.48]).

543

544    We unpacked the 3-way ANOVA into two 2-ways ANOVAs for further analysis, one for

545    each half of the experiment, with factors multisensory generalization direction (audition to

546    vision vs vision to audition) and manipulation (primary/direct vs secondary/indirect

547    modality).

548    For the first half of the experiment, our results revealed a significant main effect of

549    multisensory generalization direction ($F(1, 27)$ = 8.34, $p$ = 0.008, $\eta_p^2$ [90% CI] = 0.24 [0.04,

550    0.42]), manipulation ($F(1, 27)$ = 28.84, $p < 0.001$, $\eta_p^2$ [90% CI] = 0.52 [0.27, 0.65]) and a

551    significant interaction between multisensory generalization direction and manipulation ($F(1,$

552    27) = 12.83, $p$ = 0.001, $\eta_p^2$ [90% CI] = 0.32 [0.09, 0.50]) (left bar plot in Fig. 2D). Post-hoc

553    comparisons indicated that $\Delta RT_{Exp}$ in the primary modality of experiment 2 (i.e., vision) were

554    significantly greater than $\Delta RT_{Exp}$ in the primary modality of experiment 1 (i.e., audition)

555    ($t(27)$ = 5.72, $p < 0.001$, Cohen's $d_{av}$ [95% CI] = 1.05, [0.59, 1.50], dark blue vs dark orange

556    bars in the left bar plot of Fig. 2D), and greater than the effects of expectation in the

557    secondary modality of experiment 2 (i.e., audition) ($t(27)$ = 4.97, $p < 0.001$, Cohen's $d_{av}$

558    [95% CI] = 1.03, [0.53, 1.51], dark blue vs light orange bars in the left bar plot of Fig. 2D).

559    Moreover, the effects of expectation in the secondary modality of experiment 2 (i.e.,

560    audition) were significantly greater than those in the secondary modality of experiment 1

561    (i.e., vision) ($t(27)$ = 2.14, $p$ = 0.042, Cohen's $d_{av}$ [95% CI] = 0.48, [0.02, 0.94], light orange

562    vs light blue bars in the left bar plot of Fig. 2D).

563     For the second half of the experiment, our results only revealed a significant main effect of

564     manipulation ($F(1, 27) = 11.41$, $p = 0.002$, $\eta_p^2$ [90% CI] = 0.30 [0.07, 0.48]) (right bar plot in

565     Fig. 2D) but no significant main effect of multisensory generalization direction or significant

566     interaction between multisensory generalization direction and manipulation was found.

567

568     To summarize, these results show that: (1) an effect of generalization direction (i.e., audition

569     to vision vs vision to audition) was found only for the first half of the experiment. Here, the

570     effects of expectation generalized crossmodally in an attenuated fashion only from vision to

571     audition in experiment 2 (i.e., primary modality: vision) but no difference between audition

572     and vision was found in experiment 1 (i.e., primary modality: audition). By contrast, in the

573     second half of the experiment we did not observe an effect of multisensory generalization

574     direction. Instead, the effects of expectation generalized from the primary to the secondary

575     modality in an attenuated fashion similarly when vision or audition were the primary

576     modality. As shown in figure 2D, this difference between first and second halves can be

577     explained by the fact that observers form spatial expectations (i.e., learn signal probability

578     over space) more slowly in audition than vision. Yet, once expectations are learnt in audition,

579     the crossmodal generalization is comparable for audition and vision. In other words, the

580     effect of generalization direction that we observed in our analysis that did not yet account for

581     learning effects (i.e., our second analysis) can be explained away by the speed with which

582     observers learn signal probabilities and form spatial expectations in their primary modality.

583     In other words, signal probabilities are learnt faster in vision than audition, but once spatial

584     expectations are formed, they generalize similarly from vision to audition and vice versa.

**Figure 2**: Behavioural results of experiment 1 and 2.

Bar plots represent across subjects' mean (±SEM) RT for each of the six conditions with response requirements for experiment 1 (primary modality: audition; secondary modality: vision, **A**) and 2 (primary modality: vision; secondary modality: audition, **B**), pooling over left/right stimulus location. Overall slower RT are observed for runs type B than runs type A, which reflects differences in general response probability (see Fig. 1D) **C.** Bar plots represent across subjects' mean (±SEM) ΔRT for effects of spatial expectation (attended unexpected – attended expected hemifield) in the primary (dark bars) and secondary modalities (light bars) for experiment 1 and 2. **D.** Effects of response probability (attended unexpected – attended

595  expected hemifield) over time (i.e., first and second half: bars; consecutive sets of 2 attention

596  runs: circles) for audition and vision as primary (dark bars) or secondary (light bars)

597  modality. Brackets and stars indicate significance of main effects and interactions. *$p < 0.05$;

598  ** $p < 0.01$; *** $p < 0.001$. Audition: orange; vision: blue.

599

600  **Discussion**

601  The current study investigated how observers allocate attention and form expectations (by

602  learning signal probabilities) over space across audition and vision. We orthogonally

603  manipulated spatial attention as response requirement and expectation as stimulus probability

604  over space selectively in the primary modality and assessed their effects on behavioral

605  responses to targets presented in the primary and secondary modalities. Across two

606  experiments we alternated the assignment of vision and audition to primary or secondary

607  modality. This allowed us to compare behavioral effects of spatial attention and expectation

608  in audition and vision and their crossmodal generalization.

609  Regardless of sensory modality we observed a significant main effect of spatial attention for

610  targets in the secondary modality, in which attention was not directly manipulated. Auditory

611  spatial attention partially generalized to the visual modality and vice versa. These findings

612  converge with a large body of behavioral and neuroimaging work suggesting that attentional

613  resources are allocated interactively across the senses (Spence & Driver, 1996; Eimer &

614  Schröger, 1998; Eimer, 1999; Macaluso et al., 2002; Santangelo et al., 2009; Zuanazzi &

615  Noppeney, 2019).

616  Likewise, in the attended hemifield observers were faster at target detection in the primary

617  and secondary modalities when the hemifield was expected than unexpected (i.e., high > low

618  signal probability). Again, this response facilitation for expected (relative to unexpected)

619  spatial locations were observed irrespective of whether vision or audition served as primary

620    modality. Yet, in the unattended hemifield, in which we could assess effects of expectation

621    only for the secondary modality, we observed the opposite pattern, i.e., observers were slower

622    at target detection for expected than unexpected hemifields. Combining these two results, we

623    observed a significant interaction between spatial attention and expectation, for both vision

624    and audition as secondary modalities. In a previous study (Zuanazzi & Noppeney, 2018) we

625    argued that this interaction profile between spatial attention and expectation is explained by

626    attention and expectation jointly co-determining general response probability (i.e., the

627    probability to respond regardless of the hemifield in which the signal is presented). More

628    specifically, in runs in which attention and expectation are directed to the same hemifield

629    (runs of type A, Fig. 1A, 1C and 1D), participants have to respond to 85% of trials of the

630    entire run, but only to 65% trials in runs of type B in which attention and expectation are

631    directed to different hemifields (i.e., general response probability, Fig. 1A, 1C and 1D).

632    Hence, faster response times may result from an increase in alertness, arousal or motor

633    preparation that is needed to respond on a large proportion of trials (i.e., attended/expected

634    and unattended/unexpected conditions, run type A, Fig. 2A and 2B) (Mars et al., 2007;

635    Bestmann et al., 2008). By contrast, in runs in which attention and expectation are directed to

636    different hemifields (runs of type B, Fig. 1A, 1C and 1D), observers need to inhibit their

637    response to the frequent stimuli of the primary modality in the expected hemifield and hence

638    respond more slowly to targets in the secondary modality.

639    Critically, response probability does not depend on whether the signal is auditory or visual,

640    but it is calculated as the probability that any signal is responded to. If the expectation effects

641    result purely from amodal mechanisms (e.g., general alertness, arousal, motor preparation

642    etc.) associated with changes in response probability, the expectation effects in the attended

643    hemifield should be equal for primary and secondary modalities. By contrast, if expectations

644    (i.e., auditory or visual signal probability) are formed at least partially in a modality-specific

645    fashion, we should observe expectation effects (i.e., $\Delta RT_{Exp}$) that are greater for the primary

646    modality (where expectation was explicitly manipulated) than for the secondary modality

647    (i.e., an attenuated crossmodal generalization).

648    To arbitrate between these two hypotheses, we analyzed $\Delta RT_{Exp}$ in a 2 (multisensory

649    generalization direction) x 2 (manipulation) repeated measures ANOVA. This analysis

650    revealed a significant main effect of 'manipulation', i.e., whether the stimulus was presented

651    in the primary modality (in which expectation was explicitly manipulated) or in the

652    secondary modality. Consistent with additional modality-specific mechanisms of expectation,

653    observers showed a greater expectation effect for targets in the primary than the secondary

654    modality. These results strongly suggest that implicitly learned spatial expectations modulate

655    perceptual decision making via both modality-specific and amodal response mechanisms.

656    This duality of modality-specific and amodal mechanisms converge with recent

657    neuroimaging findings which showed effects of expectation selective for auditory stimuli as

658    primary modality in auditory cortices and higher-order frontoparietal systems (Zuanazzi &

659    Noppeney, 2019). Potentially, the activation increases in auditory cortices may reflect

660    prediction error signals based on modality-specific expectations (Friston, 2005), while higher

661    frontoparietal systems may be associated with additional response-related processes.

662    Surprisingly, the repeated measures ANOVA also revealed a main effect of multisensory

663    generalization direction with greater effects in experiment 2 (i.e., vision to audition) than

664    experiment 1 (i.e., audition to vision). Our time course analysis showed that this difference

665    between experiments does not reflect genuine differences in the effectiveness with which

666    spatial expectations are implicitly learnt and generalize from audition to vision and vice

667    versa, but reflects differences in the speed with which spatial expectations are learnt in

668    audition and vision. In the second half of the experiment in which the expectation effect in the

669    primary modality is comparable between auditory (experiment 1) and visual (experiment 2)

670     targets, we no longer observe a significant effect of multisensory generalization direction. In

671     other words, observers are slower at learning spatial expectations (i.e., signal probability) in

672     audition than in vision. But, once they have formed spatial expectations of comparable

673     precision for audition and vision, they also generalize similarly from audition to vision and

674     vice versa.

675     The difference in perceptual learning rates between vision and audition may result from how

676     the brain forms spatial representations in vision and audition (Neumann et al., 1986). While

677     visual cortices are retinotopically organized and hence directly represent spatial location in a

678     place code (i.e., based on space, Sereno et al., 1995; Maier & Groh, 2009), primary auditory

679     cortices are tonotopically organized (i.e., based on frequency, Lauter et al., 1985,

680     Middlebrooks &, 1991; Maier & Groh, 2009). In audition, spatial locations are computed

681     only indirectly from binaural amplitude and latency differences and from monoaural filtering

682     cues. Moreover, visual objects tend to be more permanent across time, whereas source sounds

683     are often transient and dynamic (Neisser, 1976; Neumann et al., 1986). Most importantly, in

684     everyday life vision provides typically (i.e., under optimal lighting conditions) more reliable

685     spatial information than audition (Dacey et al., 1992; Knudsen & Brainard, 1995; Stephen et

686     al., 2002; Talsma et al., 2008; Mengotti et al., 2018; see also Molholm et al., 2007). In the

687     current study, the high spatial reliability of the visual stimulus (i.e., a white disc) may also

688     have contributed to the shorter time for participants to learn spatial signal probabilities and

689     thus become aware of their manipulation in vision. Critically, participants' awareness of such

690     manipulation is evidenced by our questionnaires' results. Alternatively, even in absence of

691     explicit awareness, the distribution of events or targets across space could have been learnt

692     faster for more reliable visual than auditory signals (Miller & Pachella, 1973; Jabar &

693     Anderson, 2015). Conversely, in a paradigm that investigates temporal attention/expectation

694     mechanisms, the high temporal resolution and precision of auditory signals (Shimojo &

695  Shams, 2001) could facilitate learning of temporal probability in audition more than in vision

696  and crossmodal effects could change accordingly. The existence of cross-modal effects of

697  temporal attention is shown in previous studies investigating how attention is oriented to

698  different points in time (e.g., Lange & Röder, 2006). However, while the effects of spatial

699  and temporal attention were similar for auditory processing, they differed for visual and

700  tactile modalities, suggesting the existence of modality specific mechanisms also in the

701  temporal domain. To better understand the fine-grained temporal aspects of spatial

702  expectation or signal probability learning across sensory modalities, future studies will need

703  to characterize the time course of spatial learning across sensory systems.

704  So far, we have discussed that spatial expectations generalize only partially across the senses.

705  One critical question is whether this partial generalization is generic or arises because the

706  decisions rely on different processes in our paradigm. Most importantly, as we have indicated

707  in the Results section, observers had to respond to stimuli in the primary modality only in one

708  hemifield, but in the secondary modality in both hemifields. This experimental choice

709  enabled us to assess the additive and interactive effects of spatial attention and expectation in

710  both hemifields for the secondary modality. As a consequence, however, observers needed to

711  determine the hemifield in which the stimulus occurred before making a response only for the

712  primary modality. By contrast, they could respond non-discriminatively to all stimuli in the

713  secondary modality. This difference in the decision-making process most likely explains that

714  observers were faster to respond to sounds when audition was the secondary than the primary

715  sensory modality. An outstanding question is whether this difference in the decision-making

716  process can also explain the partial generalization of the expectation effects. In other words,

717  would we observe more extensive or perhaps complete generalization across sensory

718  modalities if both primary and secondary modalities rely on similar decision-making

719  processes? Or even more fundamentally, can the differences in magnitude in the expectation

720 effects for primary and secondary modality be explained by the fact that expectations

721 influence spatial discrimination processes that are required only for responding to stimuli in

722 the primary modality? To address this question, future studies may manipulate attention and

723 expectation in both sensory modalities. For instance, they may present observers with

724 auditory and visual stimuli in left and right hemifields. Auditory stimuli may occur mainly in

725 the left and visual in the right hemifield. Importantly, observers will need to respond to

726 auditory and visual stimuli only when they occur in one particular (e.g., left) hemifield, so

727 that responses to both auditory and visual stimuli will require spatial discrimination between

728 hemifields. However, as we have argued in a previous study, manipulating response

729 requirement and spatial expectations orthogonally across sensory modalities may interact at

730 the several levels by jointly specifying not only observers' general response probability but

731 also spatially selective response probabilities (Zuanazzi & Noppeney, 2018). Further, it is

732 important to emphasize that these manipulations of spatial signal probability and response

733 requirement over space operate bidirectionally from audition to vision and vice versa as well

734 as from primary to secondary modality and vice versa, thus interpretational ambiguities may

735 remain. Alternatively, complementary insights may be gained from neuroimaging research

736 that can implicitly assess the multisensory generalization of neural representations linked

737 with spatial expectations even when no response is required.

738 In summary, our results suggest that the brain allocates spatial attention and forms spatial

739 expectation to some extent interactively across audition and vision (Eimer & Schröger, 1998;

740 Eimer, 1999; Macaluso, 2010). With respect to spatial attention, our results corroborate

741 previous research. With respect to spatial expectation, we show that they rely on modality-

742 specific and amodal mechanisms. In support of modality-specific mechanisms we

743 demonstrate that spatial expectations in the attended hemifield generalize from the primary to

744 the secondary modality only in an attenuated fashion. In support of amodal response-related

33

745    mechanisms, we demonstrate that, for both primary and secondary modalities, response times

746    are closely related to the general response probability and associated processes of arousal,

747    alertness and motor preparation. Critically, our learning analysis suggests that observers learn

748    spatial probabilities more slowly in audition than vision, which may be related to their

749    different spatial reliabilities. Once observers have formed comparable spatial expectations in

750    audition, these generalize equally effectively from audition to vision as from vision to

751    audition. In other words, our results demonstrate crossmodal interactions of perceptual

752    learning (i.e., expectations building) in spatial perception but also show differences between

753    sensory modalities in terms of the speed with which signal probabilities over space are learnt.

754

755    **Declaration of Conflicting Interests**

756    The authors declared that they had no conflicts of interest with respect to their authorship or

757    the publication of this article.

758

764

765

766

767

768

769

**References**

Aller, M., Giani, A., Conrad, V., Watanabe, M., & Noppeney, U. (2015). A spatially collocated sound thrusts a flash into awareness. *Frontiers in Integrative Neuroscience*, 9(February), 1–8.

Aller, M., & Noppeney, U. (2019). To integrate or not to integrate: Temporal dynamics of hierarchical Bayesian causal inference. *PLoS biology* (Vol. 17).

Batson, M. A., Watanabe, T., Seitz, A. R., & Beer, A. L. (2011). Spatial shifts of audio-visual interactions by perceptual learning are specific to the trained orientation and eye. *Seeing and Perceiving*, 24(6), 579–594.

Bestmann, S., Harrison, L. M., Blankenburg, F., Mars, R. B., Haggard, P., Friston, K. J., & Rothwell, J. C. (2008). Influence of uncertainty and surprise on human corticospinal excitability during preparation for action. *Current Biology*, 18(10), 775–780.

Blignaut, P. (2009). Fixation identification: the optimum threshold for a dispersion algorithm. Attention, *Perception & Psychophysics*, 71(4), 881–895.

Bratzke, D., Seifried, T., & Ulrich, R. (2012). Perceptual learning in temporal discrimination: Asymmetric cross-modal transfer from audition to vision. *Experimental Brain Research*, 221(2), 205–210.

Bressler, S. L., Tang, W., Sylvester, C. M., Shulman, G. L., & Corbetta, M. (2008). Top-down control of human visual cortex by frontal and parietal cortex in anticipatory visual spatial attention. *The Journal of Neuroscience*, 28, 10056–10061.

Bueti, D., & Buonomano, D. V. (2014). Temporal perceptual learning. *Timing and Time Perception*, 2(3), 261–289.

Bueti, D., Lasaponara, S., Cercignani, M., & Macaluso, E. (2012). Learning about time: Plastic changes and interindividual brain differences. *Neuron*, 75(4), 725–737.

794    Carrasco, M. (2011). Visual attention: The past 25 years. *Vision Research*, 51(13), 1484–

795    1525.

796    Crist, R. E., Kapadia, M. K., Westheimer, G., & Gilbert, C. D. (1997). Perceptual learning of

797    spatial localization: Specificity for orientation, position, and context. *Journal of*

798    *Neurophysiology*, 78(6), 2889–2894.

799    Dacey, D. M., & Petersen, M. R. (1992). Dendritic field size and morphology of midget and

800    parasol ganglion cells of the human retina. *Proceedings of the National Academy of Sciences*

801    *of the United States of America*, 89(20), 9666–9670.

802    Doricchi, F., Macci, E., Silvetti, M., & Macaluso, E. (2010). Neural correlates of the spatial

803    and expectancy components of endogenous and stimulus-driven orienting of attention in the

804    Posner task. *Cerebral Cortex*, 20(7), 1574–1585.

805    Driver, J., & Spence, C. (1998). Attention and the crossmodal construction of space. *Trends*

806    *in Cognitive Sciences*, 2(7), 254–262.

807    Eimer, M. (1999). Can attention be directed to opposite locations in different modalities? An

808    ERP study. *Clinical Neurophysiology*, 110(7), 1252–1259.

809    Eimer, M., & Driver, J. (2001). Crossmodal links in endogenous and exogenous spatial

810    attention: evidence from event-related brain potential studies. *Neuroscience and*

811    *Biobehavioral Reviews*, 25(6), 497–511.

812    Eimer, M., & Schroger, E. (1998). ERP effects of intermodal attention and cross-modal links

813    in spatial attention. *Psychophysiology*, 35(3), 313–27.

814    Feldman, H., & Friston, K. J. (2010). Attention, uncertainty, and free-energy. *Frontiers in*

815    *Human Neuroscience*, 4(December), 1–23.

816    Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using

817    G*Power 3.1: tests for correlation and regression analyses. *Behavior Research Methods*,

818    41(4), 1149–60.

819  Friston, K. (2005). A theory of cortical responses. *Philosophical Transactions of the Royal*

820  *Society of London. Series B, Biological Sciences*, 360(1456), 815–836.

821  Gold, J. I., & Shadlen, M. N. (2001). Neural computations that underlie decisions about

822  sensory stimuli. *Trends in Cognitive Sciences*, 5(1), 10–16.

823  Greene, A. J., Easton, R. D., & LaShell, L. S. R. (2001). Visual-auditory events: Cross-modal

824  perceptual priming and recognition memory. *Consciousness and Cognition*, 10(3), 425–435.

825  Jabar, S. B., & Anderson, B. (2015). Probability shapes perceptual precision: A study in

826  orientation estimation. *Journal of Experimental Psychology: Human Perception and*

827  *Performance*, 41(6), 1666–1679.

828  Jones, S. A., Beierholm, U., Meijer, D., & Noppeney, U. (2019). Older adults sacrifice

829  response speed to preserve multisensory integration performance. *Neurobiology of Aging*, 84,

830  148–157.

831  Kim, R. S., Seitz, A. R., & Shams, L. (2008). Benefits of stimulus congruency for

832  multisensory facilitation of visual learning. *PLoS ONE*, 3(1).

833  Kincade, J. M., Abrams, R. A., Astafiev, S. V., Shulman, G. L., & Corbetta, M. (2005). An

834  event-related functional magnetic resonance imaging study of voluntary and stimulus-driven

835  orienting of attention. *The Journal of Neuroscience*, 25(18), 4593–4604.

836  Kleiner, M., Brainard, D., Pelli, D., Ingling, A., Murray, R., & Broussard, C. (2007). What's

837  new in psychtoolbox-3. *Perception*, 36(14), 1-16.

838  Knudsen, E., & Brainard, M. S. (1995). Creating a unified representation of visual and

839  auditory space in the brain. *Annual Review of Neuroscience*, 18(1), 19–43.

840  Lapid, E., Ulrich, R., & Rammsayer, T. (2009). Perceptual learning in auditory temporal

841  discrimination: No evidence for a cross-modal transfer to the visual modality. *Psychonomic*

842  *Bulletin and Review*, 16(2), 382–389.

843     Lauter, J. L., Herscovitch, P., Formby, C., & Raichle, M. E. (1985). Tonotopic organization

844     in human auditory cortex revealed by positron emission tomography. *Hearing Res*, 20, 199–

845     205.

846     Macaluso, E., Frith, C. D., & Driver, J. (2002). Supramodal effects of covert spatial orienting

847     triggered by visual or tactile events. *Journal of Cognitive Neuroscience*, 14(3), 389–401.

848     Macaluso, E. (2010). Orienting of spatial attention and the interplay between the senses.

849     *Cortex*, 46(3), 282–297.

850     Lange, K., & Röder, B. (2006). Orienting attention to points in time improves stimulus

851     processing both within and across modalities. Journal of Cognitive Neuroscience, 18(5), 715–

852     729.

853     Macaluso, E., & Doricchi, F. (2013). Attention and predictions: control of spatial attention

854     beyond the endogenous-exogenous dichotomy. *Frontiers in Human Neuroscience*,

855     7(October), 685.

856     Meegan, D. V, & Aslin, R. N. (2000). Motor timing learned. *Nature Neuroscience*, 3(9).

857     Maier, J. X., & Groh, J. M. (2009). Multisensory guidance of orienting behavior. *Hearing

858     Research*, 258(1–2), 106–112.

859     Mars, R. B., Bestmann, S., Rothwell, J. C., & Haggard, P. (2007). Effects of motor

860     preparation and spatial attention on corticospinal excitability in a delayed-response paradigm.

861     *Experimental Brain Research*, 182(1), 125–129.

862     Meijer, D., Veselič, S., Calafiore, C., & Noppeney, U. (2019). Integration of audiovisual

863     spatial signals is not consistent with maximum likelihood estimation. *Cortex*.

864     Middlebrooks, J. C., & Green, D. M. (1991). Sound localization by human listeners. *Annu.

865     Rev. Psychol.*, (42), 135–159.

866  Mengotti, P., Boers, F., Dombert, P. L., Fink, G. R., & Vossel, S. (2018). Integrating

867  modality-specific expectancies for the deployment of spatial attention. *Scientific Reports*,

868  8(1), 1210.

869  Miller, J. O., & Pachella, R. G. (1973). Locus of the stimulus probability effect. *Journal of*

870  *Experimental Psychology*, 101(2), 227–231.

871  Molholm, S., Martinez, A., Shpaner, M., & Foxe, J. J. (2007). Object-based attention is

872  multisensory: Co-activation of an object's representations in ignored sensory modalities.

873  *European Journal of Neuroscience*, 26(2), 499–509.

874  Mondor, T. A., & Amirault, K. J. (1998). Effect of same- and different-modality spatial cues

875  on auditory and visual target identification. *Journal of Experimental Psychology*, 24(3), 745–

876  755.

877  Nagarajan, S. S., Blake, D. T., Wright, B. A., Byl, N., & Merzenich, M. M. (1998). Practice-

878  related improvements in somatosensory interval discrimination are temporally specific but

879  generalize across skin location, hemisphere, and modality. *Journal of Neuroscience*, 18(4),

880  1559–1570.

881  Neisser, U. (1977). Cognition and reality. *The American Journal of Psychology*, 90(3), 541–

882  543.

883  Neumann, O., Van der Heijden, A. H., & Allport, D. A. (1986). Visual selective attention:

884  Introductory remarks. *Psychological Research*, 48(4), 185–188.

885  Nichols, T.E. & Holmes, A.P. (2002) Nonparametric permutation tests for functional

886  neuroimaging: a primer with examples. *Hum. Brain Mapp.,* 15, 1-25.

887  Odegaard, B., Wozny, D. R., & Shams, L. (2015). Biases in visual, auditory, and audiovisual

888  perception of space. *PLoS Computational Biology*, 11(12), 1–23.

889  Posner, M. I. (1980). Orienting of attention. *Q J Exp Psychol*, 32(1), 3–25.

890  Rohe, T., & Noppeney, U. (2015a). Cortical hierarchies perform Bayesian Causal Inference

891  in multisensory perception. *PLOS Biology*, 13(2), e1002073.

892  Rohe, T., & Noppeney, U. (2015b). Sensory reliability shapes perceptual inference via two

893  mechanisms. *Journal of Vision*, 15, 1–38.

894  Rohe, T., & Noppeney, U. (2016). Distinct computational principles govern multisensory

895  integration in primary sensory and association cortices. *Current Biology*, 26(4), 509–514.

896  Rohe, T., & Noppeney, U. (2018). Reliability-weighted integration of audiovisual signals can

897  be modulated by top-down control. *Eneuro*, 5, 1–20.

898  Rohenkohl, G., Gould, I. C., Pessoa, J., & Nobre, A. C. (2014). Combining spatial and

899  temporal expectations to improve visual perception. *Journal of Vision*, 14(4), 1–13.

900  Santangelo, V., Olivetti Belardinelli, M., Spence, C., & Macaluso, E. (2009). Interaction

901  between voluntary and stimulus-driven spatial attention mechanism across sensory

902  modalities. *Journal of Cognitive Neuroscience*, 21(12), 2384–2397.

903  Sereno, M. I., Dale, a M., Reppas, J. B., Kwong, K. K., Belliveau, J. W., Brady, T. J., …

904  Tootell, R. B. H. (1995). Borders of multiple visual areas in humans revealed by functional

905  magnetic resonance imaging borders of multiple visual areas in humans revealed by

906  functional magnetic resonance imaging. *Science*, 268(May), 889–893.

907  Shimojo, S., & Shams, L. (2001). Sensory modalities are not separate modalities: Plasticity

908  and interactions. *Current Opinion in Neurobiology*, 11(4), 505–509.

909  Spence, C., & Driver, J. (1996). Audiovisual links in endogenous covert spatial attention.

910  *Journal of Experimental Psychology: Human Perception and Performance*, 22(4), 1005–

911  1030.

912  Spence, C., & Driver, J. (1997). Audiovisual links in exogenous covert spatial orienting.

913  *Perception & Psychophysics*, 59(1), 1–22.

914  Stephen, J. M., Aine, C. J., Christner, R. F., Ranken, D., Huang, M., & Best, E. (2002).

915  Central versus peripheral visual field stimulation results in timing differences in dorsal stream

916  sources as measured with MEG. *Vision Research*, 42(28), 3059–3074.

917  Summerfield, C., & Egner, T. (2009). Expectation (and attention) in visual cognition. *Trends*

918  *in Cognitive Sciences*, 13(9), 403–409.

919  Talsma, D., Kok, A., Slagter, H. A., & Cipriani, G. (2008). Attentional orienting across the

920  sensory modalities. *Brain and Cognition*, 66(1), 1–10.

921  van Ede, F., de Lange, F. P., & Maris, E. (2012). Attentional cues affect accuracy and

922  reaction time via different cognitive and neural processes. *The Journal of Neuroscience*,

923  32(30), 10408–10412.

924  Vossel, S., Thiel, C. M., & Fink, G. R. (2006). Cue validity modulates the neural correlates of

925  covert endogenous orienting of attention in parietal and frontal cortex. *NeuroImage*, 32(3),

926  1257–1264.

927  Wahn, B., & König, P. (2015). Audition and vision share spatial attentional resources, yet

928  attentional load does not disrupt audiovisual integration. *Frontiers in Psychology*, 6(July), 1–

929  12.

930  Wahn, B., & König, P. (2017). Is attentional resource allocation across sensory modalities

931  task-dependent? *Advances in Cognitive Psychology*, 13(1), 83–96.

932  Ward, L. M., McDonald, J. J., & Lin, D. (2000). On asymmetries in cross-modal spatial

933  attention orienting. *Perception & Psychophysics*, 62(6), 1258–1264.

934  Warm, J. S., Stutz, R. M., & Vassolo, P. A. (1975). Intermodal transfer in temporal

935  discrimination. *Perception & Psychophysics*, 18(4), 281–286.

936  Zuanazzi, A., & Noppeney, U. (2018). Additive and interactive effects of spatial attention

937  and expectation on perceptual decisions. *Scientific Reports*, 1–12.

938    Zuanazzi, A., & Noppeney, U. (2019). Distinct neural mechanisms of spatial attention and

939    expectation guide perceptual inference in a multisensory world. *The Journal of Neuroscience*,

940    39(12), 2873–18.

941    Zuanazzi, A., & Noppeney, U. (2020). The intricate interplay of spatial attention and

942    expectation: a multisensory perspective. *Multisensory Research*, 33(4–5), 383–416.