1
2
3      Towards understanding how we pay attention in naturalistic visual search settings
4
5

6      Nora Turoman[1,2,3], Ruxandra I. Tivadar[1,4,5], Chrysa Retsa[1,6], Micah M. Murray[1,4,6,7], Pawel J.
7      Matusz[1,2,7] *
8
9
10     [1] The LINE (Laboratory for Investigative Neurophysiology), Department of Radiology,
11     Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland
12     [2] MEDGIFT Lab, Institute of Information Systems, School of Management, HES-SO Valais-
13     Wallis University of Applied Sciences and Arts Western Switzerland, Techno-Pôle 3, 3960
14     Sierre, Switzerland
15     [3] Working Memory, Cognition and Development lab, Department of Psychology and
16     Educational Sciences, University of Geneva, Geneva, Switzerland
17     [4] Department of Ophthalmology, Fondation Asile des Aveugles, Lausanne, Switzerland
18     [5] Cognitive Computational Neuroscience group, Institute of Computer Science, Faculty of
19     Science, University of Bern, Switzerland
20     [6] CIBM Center for Biomedical Imaging, Lausanne University Hospital and University of
21     Lausanne, Lausanne, Switzerland
22     [7] Department of Hearing and Speech Sciences, Vanderbilt University, Nashville, TN, USA
23
24
25
26     * Corresponding author:
27     Pawel Matusz
28     Information Systems Institute
29     University of Applied Sciences Western Switzerland (HES-SO Valais)
30     Rue Technopole 3
31     3960 Sierre, Switzerland
32     Tel: +41 27 606 9060
33
34

48
49                                        **Abstract**
50
51    Research on attentional control has largely focused on single senses and the importance of
52    behavioural goals in controlling attention. However, everyday situations are multisensory
53    and contain regularities, both likely influencing attention. We investigated how visual
54    attentional capture is simultaneously impacted by top-down goals, the multisensory nature
55    of stimuli, *and* the contextual factors of stimuli's semantic relationship and temporal
56    predictability. Participants performed a multisensory version of the Folk et al. (1992) spatial
57    cueing paradigm, searching for a target of a predefined colour (e.g. a red bar) within an
58    array preceded by a distractor. We manipulated: 1) stimuli's goal-relevance via distractor's
59    colour (matching vs. mismatching the target), 2) stimuli's multisensory nature (colour
60    distractors appearing alone vs. with tones), 3) the relationship between the distractor sound
61    and colour (arbitrary vs. semantically congruent) and 4) the temporal predictability of
62    distractor onset. Reaction-time spatial cueing served as a behavioural measure of
63    attentional selection. We also recorded 129-channel event-related potentials (ERPs),
64    analysing the distractor-elicited N2pc component both canonically and using a multivariate
65    electrical neuroimaging framework. Behaviourally, arbitrary target-matching distractors
66    captured attention more strongly than semantically congruent ones, with no evidence for
67    context modulating multisensory enhancements of capture. Notably, electrical
68    neuroimaging of surface-level EEG analyses revealed context-based influences on attention
69    to both visual and multisensory distractors, in how strongly they activated the brain and
70    type of activated brain networks. For both processes, the context-driven brain response
71    modulations occurred long before the N2pc time-window, with topographic (network-
72    based) modulations at ~30ms, followed by strength-based modulations at ~100ms post-
73    distractor onset. Our results reveal that both stimulus meaning and predictability modulate
74    attentional selection, and they interact while doing so. Meaning, in addition to temporal
75    predictability, is thus a second source of contextual information facilitating goal-directed
76    behaviour. More broadly, in everyday situations, attention is controlled by an interplay
77    between one's goals, stimuli's perceptual salience, meaning and predictability. Our study
78    calls for a revision of attentional control theories to account for the role of contextual and
79    multisensory control.
80
81    *Keywords:* attentional control, multisensory, real-world, semantic congruence, temporal
82    predictability, context

<div style="text-align:center">**Introduction**</div>

83
84
85    Goal-directed behaviour depends on the ability to allocate processing resources towards the
86    stimuli important to current behavioural goals ("attentional control"). On the one hand, our
87    current knowledge about attentional control may be limited to the rigorous, yet artificial,
88    conditions in which it is traditionally studied. On the other hand, findings from studies
89    assessing attentional control with naturalistic stimuli (audiostories, films) may be limited by
90    confounds from other processes present in such settings. Here, we systematically tested
91    how traditionally studied goal- and salience-based attentional control interact with more
92    naturalistic, context-based mechanisms.
93            In the real world, the location of goal-relevant information is rarely known in
94    advance. Since the pioneering visual search paradigm (Treisman & Gelade, 1980), we know
95    that in multi-stimulus settings target attributes can be used to control attention. Here,
96    research provided conflicting results as to whether primacy in controlling attentional
97    selection lies in task-relevance of objects' attributes (Folk et al., 1992) or their bottom-up
98    salience (e.g. Theeuwes, 1991). Folk et al. (1992) used a version of the spatial cueing
99    paradigm and revealed that attentional capture is elicited only by distractors that matched
100   the target colour. Consequently, they proposed the 'task-set contingent attentional capture'
101   hypothesis, i.e., salient objects will capture attention only if they share features with the
102   target and are thus potentially task-relevant. However, subsequently mechanisms beyond
103   goal-relevance were shown to serve as additional sources of attentional control, e.g.,
104   spatiotemporal and semantic information within the stimulus and the environment where it
105   appears (e.g., Chun & Jiang 1998; Peelen & Kastner, 2014; Summerfield et al., 2006; van
106   Moorselaar & Slagter 2019; Press et al. 2020), and multisensory processes (Matusz & Eimer,
107   2011, 2013; Matusz et al. 2015a; Lunn et al. 2019; Soto-Faraco et al. 2019).
108           Some multisensory processes occur at early latencies (<100ms post-stimulus),
109   generated within primary cortices (e.g., Talsma & Woldroff, 2005; Raij et al. 2010; Cappe et
110   al. 2010; reviewed in de Meo et al., 2015; Murray et al. 2016a). This enables multisensory
111   processes to influence attentional selection in a bottom-up fashion, potentially
112   independently of the observer's goals. This idea was supported by Matusz and Eimer (2011)
113   who used a multisensory adaptation of Folk et al.'s (1992) task. The authors replicated the
114   task-set contingent attentional capture effect and showed that visual distractors captured
115   attention more strongly when accompanied by a sound, regardless of their goal-relevance.
116   This demonstrated the importance of bottom-up multisensory enhancement for attentional
117   selection of visual objects. However, interactions between such goals, multisensory
118   influences on attentional control, and the stimuli's temporal and semantic context[1] remain
119   unknown.
120
121   **Top-down contextual factors in attentional control**
122   The temporal structure of the environment is routinely used by the brain to build
123   predictions. Attentional control uses such predictions to improve the selection of target
124   stimuli (e.g., Correa et al., 2005; Coull et al., 2000; Green & McDonald, 2010; Miniussi et al.,
125   1999; Naccache et al., 2002; Rohenkohl et al., 2014; Tivadar et al. 2021) and the inhibition of

---

[1] Context has been previously defined as the "immediate situation in which the brain operates… shaped by external circumstances, such as properties of sensory events, and internal factors, such as behavioural goal, motor plan, and past experiences" (van Atteveldt et al., 2014).

126    task-irrelevant stimuli (here, location- and feature-based predictions have been more
127    researched than temporal predictions; e.g., reviewed in Noonan et al. 2018; van Moorselaar
128    & Slagter 2020a). In naturalistic, multisensory settings, temporal predictions are known to
129    improve language comprehension (e.g. Luo & Poeppel, 2007; ten Oever & Sack, 2015), yet
130    their role as a source of attentional control is less known (albeit see, Zion Golumbic et al.
131    2012, for their role in the "cocktail party" effect). Semantic relationships are another basic
132    principle of organising information in real-world contexts. Compared to semantically
133    incongruent or meaningless (arbitrary) multisensory stimuli, semantically congruent stimuli
134    are more easily identified and remembered (e.g. Laurienti et al. 2004; Murray et al., 2004;
135    Doehrmann & Naumer 2008; Chen & Spence, 2010; Matusz et al., 2015a; Tovar et al. 2020;
136    reviewed in ten Oever et al. 2016; Murray et al., 2016b; Matusz et al. 2020) and also, more
137    strongly attended (Matusz et al. 2015b, 2019a, 2019b; reviewed in Soto-Faraco et al., 2019;
138    Matusz et al. 2019c). For example, Iordanenscu et al. (2009) demonstrated that search for
139    naturalistic objects is faster when accompanied by irrelevant albeit congruent sounds.
140        What is unclear from existing research is the degree to which goal-based attentional
141    control interacts with salience-driven (multisensory) mechanisms *and* such contextual
142    factors. Researchers have been clarifying such interactions, but typically in a pair-wise
143    fashion, between e.g., attention and semantic memory, or attention and predictions
144    (reviewed in Summerfield & Egner 2009; Nobre & Gazzaley 2016; Press et al. 2020).
145    However, in everyday situations these processes do not interact in an orthogonal, but,
146    rather, a synergistic fashion, with multiple sources of control interacting simultaneously (ten
147    Oever et al. 2016; Nastase et al. 2020). Additionally, in the real world, these processes
148    operate on both unisensory and multisensory stimuli, where the latter are often more
149    perceptually salient than the former (e.g., Santangelo & Spence 2007; Matusz & Eimer
150    2011). Thus, one way to create more complete and "naturalistic" theories of attentional
151    control is by investigating how one's goals interact with *multiple* contextual factors in
152    controlling attentional selection – and doing so in *multi-sensory* settings.
153
154    **The present study**
155    To shed light on how attentional control operates in naturalistic visual search settings, we
156    investigated how visual and multisensory attentional control interact with distractor
157    temporal predictability and multisensory semantic relationship when all are manipulated
158    simultaneously. We likewise set out to identify brain mechanisms supporting such complex
159    interactions. To address these questions in a rigorous and state-of-the-art fashion, we
160    employed a 'naturalistic laboratory' approach that builds on several methodological
161    advances (Matusz et al., 2019c). First, we used a paradigm that isolates a specific cognitive
162    process, i.e., Matusz and Eimer's (2011) multisensory adaptation of the Folk et al.'s (1992)
163    task, where we additionally manipulated distractors' temporal predictability and
164    relationship between their auditory and visual features. In Folk et al.'s task, attentional
165    control is measured via well-understood spatial cueing effects, where larger effects (e.g., for
166    target-colour and audiovisual distractors) reflect stronger attentional capture. Notably,
167    distractor-related responses have the added value as they isolate attentional from later,
168    motor response-related, processes. Second, we measured a well-researched brain correlate
169    of attentional object selection, the N2pc event-related potential (ERP) component. The
170    N2pc is a negative-going voltage deflection starting at around 200ms post-stimulus onset at
171    posterior electrode sites contralateral to stimulus location (Luck & Hillyard, 1994a, 1994b;
172    Eimer, 1996; Girelli & Luck, 1997). Studies canonically analysing N2pc have provided strong

173    evidence for task-set contingence of attentional capture (e.g., Kiss et al., 2008a; 2008b;
174    Eimer et al., 2009). Importantly, N2pc is also sensitive to meaning (e.g., Wu et al., 2015) and
175    predictions (e.g., Burra & Kerzel, 2013), whereas its sensitivity to multisensory enhancement
176    is limited (van der Burg et al. 2011, but see below). This joint evidence makes the N2pc a
177    valuable 'starting point' for investigating interactions between visual goals and more
178    naturalistic sources of control. Third, analysing the traditional EEG markers of attention with
179    advanced frameworks like electrical neuroimaging (e.g., Lehmann & Skrandies 1980; Murray
180    et al., 2008; Tivadar & Murray 2019) might offer an especially robust, accurate and
181    informative approach.
182        Briefly, an electrical neuroimaging framework encompasses multivariate, reference-
183    independent analyses of global features of the electric scalp field. Its main added value is
184    that it readily distinguishes the neurophysiological mechanisms driving differences in ERPs
185    across experimental conditions in *surface-level* EEG: 1) "gain control" mechanisms,
186    modulating the strength of activation within an indistinguishable brain network, and 2)
187    topographic (network-based) mechanisms, modulating the recruited brain sources (scalp
188    EEG topography differences forcibly follow from changes in the underlying sources; Murray
189    et al. 2008). Electrical neuroimaging overcomes interpretational limitations of canonical
190    N2pc analyses. Most notably, a difference in mean N2pc amplitude can arise from both
191    strength-based and **topographic** mechanisms (albeit it is assumed to signify gain control); it
192    can also emerge from different brain source configurations (for a full discussion, see Matusz
193    et al., 2019b).
194        We recently used this approach to better understand brain and cognitive
195    mechanisms of attentional control. We revealed that distinct brain networks are active
196    during ~N2pc time-window during visual goal-based *and* multisensory bottom-up attention
197    control (across the lifespan; Turoman et al. 2021a, 2021b). However, these reflect spatially-
198    selective, lateralised brain mechanisms, partly captured by the N2pc (via the contra- and
199    ipsilateral comparison). There is little existing evidence to strongly predict how interactions
200    between goals, stimulus salience and context can occur in the brain. Schröger et al. (2015)
201    proposed that temporally unpredictable events attract attention more strongly (to serve as
202    a signal to reconfigure the predictive model about the world), visible in larger behavioural
203    responses and ERP amplitudes. Both predictions and semantic memory could be used to
204    reduce attention to known (i.e., less informative) stimuli. Indeed, goal-based control uses
205    knowledge to facilitate visual and multisensory processing (Summerfield et al. 2008;
206    Iordanescu et al., 2008; Matusz et al. 2016; Sarmiento et al. 2016). However, several
207    questions remain. Does knowledge affect attention to task-*irrelevant* stimuli the same way?
208    How early do contextual factors influence stimulus processing here, if both processes are
209    known to do so <150ms post-stimulus (Summerfield & Egner, 2009; ten Oever et al. 2016).
210    Finally, do contextual processes operate through lateralised or non-lateralised brain
211    mechanisms? Below we specify our hypotheses.
212        We expected to replicate the TAC[2] effect: In behaviour, visible as large behavioural
213    capture for target-colour matching distractors and no capture for nontarget-colour
214    matching distractors (e.g., Folk et al., 1992; Folk, et al., 2002; Lien et al., 2008); in canonical
215    EEG analyses - enhanced N2pc amplitudes for target-colour than nontarget-colour
216    distractors (Eimer et al., 2009). TAC should be modulated by both contextual factors: the

---

[2] Please see Appendix 1 for the full list of abbreviations used in the manuscript.

217    predictability of distractor onset and the multisensory relationship between distractor
218    features (semantic congruence vs. arbitrary pairing; Wu et al. 2015; Burra & Kerzel, 2013).
219    However, as discussed above, we had no strong predictions how the contextual factors
220    would modulate TAC (or if they interact while doing so), as these effects have never been
221    tested systematically together, on audio-visual and task-irrelevant stimuli. For multisensory
222    enhancement of capture, we expected to replicate it behaviourally (Matusz & Eimer 2011),
223    but without strong predictions about concomitant N2pc modulations (c.f. van der Burg et al.
224    2011). We expected multisensory enhancement of capture to be modulated by contextual
225    factors, especially multisensory relationship, based on the extensive literature on the role of
226    semantic congruence in multisensory cognition (Doehrmann & Naumer, 2008; ten Oever et
227    al. 2016). Again, we had no strong predictions as to the directionality of these modulations
228    or interaction of their influences.
229        We were primarily interested if interactions between visual goals (task-set
230    contingent attentional capture, TAC), multisensory salience (multisensory enhancement of
231    capture, MSE) and contextual processes are supported by strength-based (i.e., "gain"-like;
232    i.e., one network is active more strongly for some and less strongly for other experimental
233    conditions) and/or topographic (i.e., different networks are activated for different
234    experimental conditions) brain mechanisms, as observable in *surface-level* EEG data when
235    using multivariate analyses like electrical neuroimaging. The second aim of our study was to
236    clarify if the attentional and contextual control interactions are supported by lateralised
237    (N2pc-like) or nonlateralized mechanisms. To this aim, we analysed if those interactions are
238    captured by canonical N2pc analyses or electrical neuroimaging analyses of the lateralised
239    distractor-elicited ERPs ~180-300ms post-stimulus (N2pc-like time-window). These analyses
240    would reveal presence of strength- and topographic *spatially-selective* brain mechanisms
241    contributing to attentional control. However, analyses of the N2pc assume not only
242    lateralised activity, but also symmetry; in brain anatomy but also in scalp electrodes,
243    detecting homologous brain activity over both hemispheres. This may prevent them from
244    detecting other, less-strongly-lateralised brain mechanisms of attentional control. We have
245    previously found nonlateralised mechanisms to play a role in attentional control in
246    multisensory settings (Matusz et al. 2019b). Also, semantic information and temporal
247    expectations (and feature-based attention) are known to modulate nonlateralised ERPs
248    (Saenz et al. 2003; Dell'Acqua et al. 2010; Dassanayake et al. 2016). Thus, as the third aim of
249    our study, we investigated whether contextual control affects stages associated with
250    attentional selection (reflected by the N2pc) or also earlier processing stages. We tested this
251    by measuring strength- and/or topographic nonlateralised brain mechanisms across the
252    whole post-stimulus time-period activity.
253
254
255                          **Materials and Methods**
256
257    **Participants**
258    Thirty-nine adult volunteers participated in the study (5 left-handed, 14 males, $M_{age}$: 27.5
259    years, *SD:* 4 years, range: 22–38 years). We conducted post-hoc power analyses for the two
260    effects that have been previously behaviourally studied with the present paradigm, namely
261    TAC and MSE. Based on the effect sizes in the original Matusz and Eimer (2011, Exp.2), the
262    analyses revealed sufficient statistical power for both behavioural effects with the collected
263    sample. For ERP analyses, we could calculate power analyses only for the TAC effect. Based

264 on a purely visual ERP study (Eimer et al., 2009) we revealed there to be sufficient statistical
265 power to detect TAC in the N2pc in the current study (all power calculations are available in
266 the Supplemental Online Materials, SOMs). Participants had normal or corrected-to-normal
267 vision and normal hearing and reported no prior or current neurological or psychiatric
268 disorders. Participants provided informed consent before the start of the testing session. All
269 research procedures were approved by the Cantonal Commission for the Ethics of Human
270 Research (CER-VD; no. 2018-00241).
271
272 **Task properties and procedures**
273 *General task procedures.* The full experimental session consisted of participants completing
274 four experimental Tasks. All the Tasks were close adaptations of the original paradigm of
275 Matusz and Eimer (2011 Exp.2; that is, in turn, an adaptation of the spatial-cueing task of
276 Folk et al. [1992]). Across all the Tasks, the instructions and the overall experimental set up
277 were the same as in the study of Matusz & Eimer (1992, Exp.2; see Figure 1A). Namely,
278 participants searched for a target of a predefined colour (e.g., a red bar) in a 4-element
279 array, and assessed the target's orientation (vertical vs. horizontal). Furthermore, in all
280 Tasks, the search array was always preceded by an array containing colour distractors.
281 Those distractors always either matched the target colour (red set of dots) or matched
282 another, nontarget colour (blue set of dots); on 50% of all trials the colour distractors would
283 be accompanied by a sound (audiovisual distractor condition). The distractor appeared in
284 each of the four stimulus locations with equal probability (25%) and was thus not predictive
285 of the location of the incoming target. Differences in response speed on trials where
286 distractor and target appeared in the same vs. different locations were used to calculate
287 behavioural cueing effects that were the basis of our analyses (see below). Like in the
288 Matusz and Eimer (2011) study, across all Tasks, each trial consisted of the following
289 sequence of arrays: base array (duration manipulated; see below), followed by distractor
290 array (50ms duration), followed by a fixation point (150ms duration), and finally the target
291 array (50ms duration, see Figure 1A).
292 　　　The differences to the original study involved the changes necessary to implement
293 the two new, contextual factors that were manipulated across the four Tasks (Figure 1B).[3]
294 To implement the *Multisensory Relationship* factor, after the first two Tasks, participants
295 completed a training session (henceforth *Training*), after which they completed the
296 remaining two Tasks. To implement the *Distractor Onset* factor, the predictability of the
297 onset of the distractors was manipulated, being either stable (as in the original study, Tasks
298 2 and 4) or varying between three durations (Tasks 1 and 3). The setup involving 4
299 consecutive Tasks separated by Training allowed a systematic comparison between the four
300 levels of the two contextual factors. We now describe in more detail the procedures related
301 to all Tasks, after which we provide more details on the different tasks themselves.
302 　　　The base array contained four differently coloured sets of closely aligned dots, each
303 dot subtending 0.1° × 0.1° of visual angle. The sets of dots were spread equidistally along
304 the circumference of an imaginary circle against a black background, at an angular distance
305 of 2.1° from a central fixation point. Each set could be of one of four possible colours

---

[3] Compared to the original paradigm, we made two additional changes, to enable the Task 1 to serve as an adult control study in a developmental study (Turoman et al., 2021). We reduced the number of elements in all arrays from 6 to 4, and targets were reshaped to look like diamonds rather than rectangles. Notably, despite these changes, we have replicated here the visual and multisensory attentional control effects.

306  (according to the RGB scale): green (0/179/0), pink (168/51/166), gold (150/134/10), silver
307  (136/136/132). In the distractor array, one of the base array elements changed colour to
308  either a target-matching colour, or a target-nonmatching colour that was not present in any
309  of the elements before. The remaining three distractor array elements did not change their
310  colour. The distractors and the subsequent target diamonds could have either a blue (RGB
311  values: 31/118/220) or red (RGB values: 224/71/52) colour. The target array contained four
312  bars (rectangles), where one was always the colour-defined target. The target colour was
313  counterbalanced across participants. Target orientation (horizontal or vertical) was
314  randomly determined on each trial. The two distractor colours were randomly selected with
315  equal probability before each trial, and the location of the colour change distractor was not
316  spatially predictive of the subsequent target location (distractor and target location were
317  the same on 25% of trials). On half of all trials, distractor onset coincided with the onset of a
318  pure sine-wave tone, presented from two loudspeakers on the left and right sides of the
319  monitor. Sound intensity was 80 dB SPL (as in Matusz & Eimer, 2011), measured using an
320  audiometer placed at a position adjacent to participants' ears (CESVA SC160). Through
321  manipulations of the in-/congruence between distractor and target colour and of the
322  presence/absence of sound during distractor presentations, there were four types of
323  distractors, across all the Tasks: visual distractors that matched the target colour (TCCV,
324  short for *target-colour cue visual*), visual distractors that did not match the target colour
325  (NCCV, *nontarget-colour cue visual*), audiovisual distractors that matched the target colour
326  (TCCAV, *target-colour cue audiovisual*), and audiovisual distractors that did not match the
327  target colour (NCCAV, *nontarget-colour cue, audiovisual*).
328      The experimental session consisted of 4 Tasks, each spanning 8 blocks of 64 trials.
329  This resulted in 2,048 trials in total (512 trials per Task). Participants were told to respond as
330  quickly and accurately as possible to the targets' orientation by pressing one of two
331  horizontally aligned round buttons (Lib Switch, Liberator Ltd.) that were fixed onto a tray
332  bag on the participants' lap. If participants did not respond within 5000ms of the target
333  onset, next trial was initiated; otherwise the next trial was initiated immediately after the
334  button press. Feedback on accuracy was given after each block, followed by a progress
335  screen (*a treasure map*), which informed participants of the number of remaining blocks
336  and during which participants could take a break. Breaks were also taken between each
337  Task, and before and after the Training. As a pilot study revealed sufficient proficiency at
338  conducting the tasks after a few trials (over 50% accuracy), participants did not practice
339  doing the Tasks before administration unless they had trouble following the task
340  instructions. The experimental session took place in a dimly lit, sound-attenuated room,
341  with participants seated at 90cm from a 23" LCD monitor with a resolution of 1080 × 1024
342  (60-Hz refresh rate, HP EliteDisplay E232). All visual elements were approximately
343  equiluminant (~20cd/m$^2$), as determined by a luxmeter placed at a position close to the
344  screen, measuring the luminance of the screen filled with each respective element's colour.
345  The averages of three measurement values per colour were averaged across colours and
346  transformed from lux to cd/m$^2$ to facilitate comparison with the results of Matusz & Eimer
347  (2011). The experimental session lasted <3h in total, including an initial explanation and
348  obtaining consent, EEG setup, administration of Tasks and Training, and breaks.
349      We now describe the details of the Tasks and Training, which occurred always in the
350  same general order: Tasks 1 and 2, followed by the Training, followed by Tasks 3 and 4 (the
351  order of Tasks 1 and 2 and, separately, the order of Tasks 3 and 4, was counterbalanced
352  across participants). Differences across the four Tasks served to manipulate the two

353    contextual factors (illustrated in Figure 1B). The factor *Multisensory Relationship*
354    represented the relation between the visual (the colour of the distractor) and the auditory
355    (the accompanying sound) component stimuli that made up the distractors. These two
356    stimuli could be related just by their simultaneous presentation (Arbitrary condition) or by
357    additionally sharing meaning (Congruent condition). The factor *Distractor Onset*
358    represented the temporal predictability of the distractors, i.e., whether their onset was
359    constant within Tasks and, therefore Predictable condition, or variable and, therefore,
360    Unpredictable condition. The manipulation of the two context factors led to the creation of
361    four contexts, represented by each of the Tasks 1 – 4 (i.e., Arbitrary Unpredictable, Arbitrary
362    Predictable, Congruent Unpredictable, and Congruent Predictable). To summarise, the two
363    within-task factors encompassing distractor colour and tone presence/absence, together
364    with the two between-task factors resulted in a total of four factors in our analysis design:
365    Distractor Colour (TCC vs. NCC), Distractor Modality (V vs. AV), Distractor Onset (Predictable
366    vs. Unpredictable) and Multisensory Relationship (Arbitrary vs. Congruent)[4].
367

368    **Tasks 1 and 2.** As mentioned above, across Tasks 1 and 2, the colour of the distractor
369    and the sound accompanying the colour distractor were related only by their simultaneous
370    presentation. As such, trials from Tasks 1 and 2 made up the Arbitrary condition of the
371    Multisensory Relationship factor. Sound frequency was always 2000Hz (as in Matusz &
372    Eimer, 2011). The main difference between Task 1 and Task 2 lied in the onset of the
373    distractors in those tasks. Unbeknownst to participants, in Task 1, duration of the base array
374    varied randomly on a trial-by-trial basis, between 100ms, 250ms and 450ms, i.e., the
375    distractor onset was unpredictable. In contrast, in Task 2, the base array duration was
376    always constant, at 450ms, i.e., the distractor onset was predictable. With this
377    manipulation, considering the between-task factors: Task 1 represented Arbitrary
378    (Multisensory Relationship) and Unpredictable (Distractor Onset) trials, and Task 2 -
379    Arbitrary (Multisensory Relationship) and Predictable (Distractor Onset) trials.
380    **Training.** The Training served to induce in participants a semantic-level association
381    between a specific distractor colour and a specific pitch. This rendered distractors in the
382    Tasks following the Training semantically related (Congruent), and distractors in the
383    preceding Tasks semantically unrelated (Arbitrary). The Training consisted of an Association
384    phase followed by a Testing phase (both based on the association task in Sui, He &
385    Humphreys, 2012; see also Sun et al., 2016).
386    *I. Association phase.* The Association phase served to induce the AV associations in
387    participants. Participants were shown alternating colour word–pitch pairs, presented in the
388    centre of the screen (the tone was presented from two lateral speakers, rendering it
389    spatially diffuse and so appearing to also come from the centre of the screen). The words
390    denoted one of two distractor colours (*red* or *blue*). The tone of either high (4000Hz) or low
391    (300Hz) pitch. Both the colour word and sound were presented for 2 seconds, after which a
392    central fixation cross was presented for 150ms, followed by the next colour word–pitch pair.
393    There could be two possible colour–pitch pairing options. In one, the high-pitch tone was

---

[4] As part of our stimulus design and alike Matusz and Eimer (2011), we manipulated a third within-task factor, i.e., whether the distractor and the upcoming target appeared in the same compared to a different location. This manipulation was necessary for us to compute behavioural attentional capture that were the bases of our complex 4-factor analyses However, to avoid confusing the reader, we have removed the descriptions of this factor from the main text and we only refer briefly to the manipulation in the *General task procedures*.

394 associated with the word *red*, the low-pitch tone - with the word *blue*. In the second option,
395 the high-pitch tone was associated with the word *blue*, the low-pitch tone with the word *red*
396 (see Figure 1C, Association phase). Pairing options were counterbalanced across
397 participants. Thus, for participants trained with the first option, the pairing of word *red* and
398 a high-pitch tone would be followed by the pairing of the word *blue* with a low-pitch tone,
399 again followed by the *red*–high pitch pairing, etc. There were 10 presentations per pair,
400 resulting in a total of 20 trials. Colour words were chosen instead of actual colours to ensure
401 that the AV associations were based on meaning rather than lower-level stimulus features
402 (for examples of such taught crossmodal correspondences see, e.g., Ernst, 2007). Also,
403 colour words were shown in participants' native language (speakers: 19 French, 8 Italian, 5
404 German, 4 Spanish, 3 English). Participants were instructed to try to memorise the pairings
405 as best as they could, being informed that they would be subsequently tested on how well
406 they learnt the pairings.
407
408
409            ** FIGURE 1 HERE **
410
411
412   *II. Testing phase.* The Testing phase served to ensure that the induced colour–pitch
413 associations was strong. Now, participants were shown colour word–pitch pairings (as in the
414 Association phase) but also colour–pitch pairings (a string of x's in either red or blue, paired
415 with a sound, Figure 1C, *Testing phase* panel). Additionally, now, the pairings either
416 matched or mismatched the type of associations induced in the Association phase, e.g., if
417 the word *red* have been paired with a high-pitch tone in the Testing phase, the matching
418 pair now would be a word *red* or red x's, paired with a high-pitch tone, and mismatching
419 pair - the word *red* or red x's paired with a low-pitch tone. Participants had to indicate if a
420 given pair was matched or mismatched by pressing one of two buttons (same button setup
421 as in the Tasks). Participants whose accuracy was ≤50% had to repeat the testing.
422   The paradigm that Sui et al. (2012) have designed led to people being able to
423 reliably associate low-level visual features (colours, geometric shapes) with abstract social
424 concepts (themselves, their friend, a stranger). Following their design, in the Testing phase,
425 each pairing was shown for 250ms, of which 50ms was the sound (instead of the stimulus
426 duration of 100ms that Sui et al. used, to fit our stimulus parameters). The pairing
427 presentation was followed by a blank screen (800ms), during which participants had to
428 respond, and after each responses a screen with feedback on their performance appeared.
429 Before each trial, a fixation cross was also shown, for 500ms. Each participant performed
430 three blocks of 80 trials, with 60 trials per possible combination (colour word – sound
431 matching, colour word – sound mismatching, colour – sound matching, colour – sound
432 mismatching). A final summary of correct, incorrect, and missed trials was shown at the end
433 of Testing phase.
434   **Tasks 3 and 4.** Following the Training, in Tasks 3 and 4, the distractors' colour and
435 the accompanying sound were now semantically related. Thus, the trials from these two
436 Tasks made up the (semantically) Congruent condition of the Multisensory Relationship
437 factor. Only congruent colour–pitch distractor pairings were now presented, as per the
438 pairing option induced in the participants. That is, if the colour red was paired with a high-
439 pitch tone in the Association phase, red AV distractors in Tasks 3 and 4 were always
440 accompanied by a high-pitch tone. The pitch of sounds was now either 300Hz (low-pitch

441 condition; chosen based on Matusz & Eimer, 2013, where two distinct sounds were used) or
442 4000Hz (high-pitch condition; chosen for its comparable perceived loudness in relation to
443 the above two sound frequencies, as per the revised ISO 226:2003 equal-loudness-level
444 contours standard; Spierer et al. 2013). As between Tasks 1 and 2, Task 3 and Task 4 differed
445 in the predictability of distractor onsets, i.e., in Task 3, distractor onset was unpredictable,
446 and in Task 4 - predictable. Therefore, Task 3 represented Congruent (Multisensory
447 Relationship) and Unpredictable (Distractor Onset) trials, and Task 4 - Congruent
448 (Multisensory Relationship) and Predictable (Distractor Onset) trials.
449
450 **EEG acquisition and preprocessing**
451 Continuous EEG data sampled at 1000Hz was recorded using a 129-channel HydroCel
452 Geodesic Sensor Net connected to a NetStation amplifier (Net Amps 400; Electrical
453 Geodesics Inc., Eugene, OR, USA). Electrode impedances were kept below 50kΩ, and
454 electrodes were referenced online to Cz. First, offline filtering involved a 0.1Hz high-pass
455 and 40Hz low-pass as well as 50Hz notch (all filters were second-order Butterworth filters
456 with −12dB/octave roll-off, computed linearly with forward and backward passes to
457 eliminate phase-shift). Next, the EEG was segmented into peri-stimulus epochs from 100ms
458 before distractor onset to 500ms after distractor onset. An automatic artefact rejection
459 criterion of ±100μV was used, along with visual inspection. Epochs were then screened for
460 transient noise, eye movements, and muscle artefacts using a semi-automated artefact
461 rejection procedure. Data from artefact contaminated electrodes were interpolated using
462 three-dimensional splines (Perrin et al., 1987). Across all Task, 11% of epochs were removed
463 on average and 8 electrodes were interpolated per participant (6% of the total electrode
464 montage).
465  Cleaned epochs were averaged, baseline corrected to the 100ms pre-distractor time
466 interval, and re-referenced to the average reference. Next, to eliminate residual
467 environmental noise in the data, a 50Hz filter was applied[5]. All the above steps were done
468 separately for ERPs from the four distractor conditions, and separately for distractors in the
469 left and right hemifield. We next relabeled ERPs from certain conditions, as is done in
470 traditional lateralised ERP analyses (like those of the N2pc). Namely, we relabelled single-
471 trial data from all conditions where distractors appeared on the *left* so that the electrodes
472 over the left hemiscalp now represented the activity over the right hemiscalp, and
473 electrodes over the right hemiscalp – represented activity over the left hemiscalp, thus
474 creating "mirror distractor-on-the-right" single-trial data. Next, these mirrored data and the
475 veridical "distractor-on-the-right" data from each of the 4 distractor conditions were
476 averaged together, creating a single average ERP for each of the 4 distractor conditions. The
477 contralaterality factor (i.e. contralateral vs. ipsilateral potentials) is normally represented by
478 separate ERPs (one for contralateral activity, and one for ipsilateral activity; logically more
479 pairs for pair-wise N2pc analyses). In our procedure, the lateralised voltage gradients across
480 the whole scalp are preserved within each averaged ERP by simultaneous inclusion of both
481 contralateral and ipsilateral hemiscalp activation. Such a procedure enabled us to fully
482 utilise the capability of the electrical neuroimaging analyses in revealing both lateralised and
483 non-lateralised mechanisms that support the interactions of attentional control with

---

[5] While filtering following epoch creation is normally discouraged (e.g., Widmann et al. 2015), control analyses we have
carried out demonstrated that our filtering procedure was necessary and did not harm the data quality within our time-
window of interest (for results of control analyses, see SOMs: Justification of filtering choices).

484    context control. As a result of the relabelling, we obtained 4 different ERPs: TCCV (target
485    colour-cue, Visual), NCCV (nontarget colour-cue, Visual), TCCAV (target colour-cue,
486    AudioVisual), NCCAV (nontarget colour-cue, AudioVisual). Preprocessing and EEG analyses,
487    unless otherwise stated, were conducted using CarTool software (available for free at
488    www.fbmlab.com/cartool-software/; Brunet, Murray, & Michel, 2011).
489
490    **Data analysis design**
491        **Behavioural analyses**. Like in Matusz and Eimer (2011), and because mean reaction
492    times (RTs) and accuracy did not differ significantly between the four Tasks, the basis of our
493    analyses was RT spatial cueing effects (henceforth "behavioural capture effects"). These
494    were calculated by subtracting the mean RTs for trials where the distractor and target were
495    in the same location from the mean RTs for trials where the distractor and the target
496    location differed, separately for each of the four distractor conditions. Such spatial cueing
497    data were analysed using the repeated-measures analysis of variance (rmANOVA). Error
498    rates (%) were also analysed. As they were not normally distributed, we analysed error rates
499    using the Kruskal–Wallis $H$ test and the Durbin test. The former was used to analyse if error
500    rates differed significantly between Tasks, while the latter was used to analyse differences
501    between experimental conditions within each Task separately.
502        Following Matusz and Eimer (2011), RT data were cleaned by discarding incorrect
503    and missed trials, as well as RTs below 200ms and above 1000ms. Additionally, to enable
504    more direct comparisons with the developmental study for which current Task 1 served as
505    an adult control (Turoman et al., 2021a, 2021b), we have further removed trials with RTs
506    outside 2.5SD of the individual mean RT. As a result, a total of 5% of trials across all Tasks
507    were removed. Next, behavioural capture effects were submitted to a four-way 2 × 2 × 2 × 2
508    rmANOVA with factors: Distractor Colour (TCC vs. NCC), Distractor Modality (V vs. AV),
509    Multisensory Relationship (Multisensory Relationship; Arbitrary vs. Congruent), and
510    Distractor Onset (Distractor Onset; Unpredictable vs. Predictable). Due to the error data not
511    fulfilling criteria for normality, we used Distractor-Target location as a factor in the analysis,
512    conducting 3-way Durbin tests for each Task, with factors Distractor Colour, Distractor
513    Modality, and Distractor-Target Location. All analyses, including post-hoc paired $t$-tests,
514    were conducted using SPSS for Macintosh 26.0 (Armonk, New York: IBM Corporation). For
515    brevity, we only present the RT results in the Results, and the error rate results can be found
516    in SOMs.
517        **ERP analyses**. The preprocessing of the ERPs triggered by the visual and audiovisual
518    distractors across the 4 different experimental blocks created ERP averages in which the
519    contralateral versus ipsilateral ERP voltage gradients across the whole scalp were preserved.
520    We first conducted a canonical N2pc analysis, as the N2pc is a well-studied and well-
521    understood correlate of attentional selection in visual settings. However, it is unclear if the
522    N2pc also indexes bottom-up attentional selection modulations by multisensory stimuli, or
523    top-down modulations by contextual factors like multisensory semantic relationships (for
524    visual-only study, see e.g., Wu et al. 2015) or stimulus onset predictability (for visual-only
525    study, see e.g., Burra & Kerzel, 2013). N2pc analyses served also to bridge electrical
526    neuroimaging analyses with the existing literature and EEG approaches more commonly
527    used to investigate attentional control. Briefly, electrical neuroimaging encompasses a set of
528    multivariate, reference-independent analyses of global features of the electric field
529    measured at the scalp (König et al., 2014; Michel & Murray, 2012; Murray, Brunet, & Michel,
530    2008; Lehmann & Skrandies, 1980; Tivadar & Murray, 2019; Tzovara et al., 2012) that can

12

531    detect spatiotemporal patterns in EEG across different contexts and populations (e.g., Neel
532    et al. 2019; Matusz et al. 2018). The key advantages of electrical neuroimaging analyses over
533    canonical N2pc analyses and how the former can complement the latter when combined,
534    are described in the Introduction.
535         *Canonical N2pc analysis.* To analyse lateralised mechanisms using the traditional
536    N2pc approach, we extracted mean amplitude values from, first, two electrode clusters
537    comprising PO7/8 electrode equivalents (e65/90; most frequent electrode pair used to
538    analyse the N2pc), and, second, their six immediate surrounding neighbours (e58/e96,
539    e59/e91, e64/e95, e66/e84, e69/e89, e70/e83), over the 180–300ms post-distractor time-
540    window (based on time-windows commonly used in traditional N2pc studies, e.g., Luck &
541    Hillyard, 1994b; Eimer, 1996; including distractor-locked N2pc, Eimer & Kiss 2008; Eimer  et
542    al. 2009). Analyses were conducted on the mean amplitude of the N2pc difference
543    waveforms, which were obtained by subtracting the average of amplitudes in the ipsilateral
544    posterior-occipital cluster from the average of amplitudes in the contralateral posterior-
545    occipital cluster. This step helped mitigate the loss of statistical power that could result from
546    the addition of contextual factors into the design. N2pc means were thus submitted to a 4-
547    way $2 \times 2 \times 2 \times 2$ rmANOVA with factors Distractor Colour (TCC vs. NCC), Distractor Modality
548    (V vs. AV), Multisensory Relationship (Arbitrary vs. Congruent), and Distractor Onset
549    (Unpredictable vs. Predictable), analogously to the behavioural analysis. Notably, the N2pc
550    is not sensitive to the location of the stimulus of interest *per se,* but rather to the side of its
551    presentation. As such, in canonical analyses of distractor-elicited N2pc, the congruence
552    between distractor and target, unlike in behavioural analyses, is not considered (e.g., Lien et
553    al. 2008; Eimer & Kiss 2008; Eimer et al. 2009). Consequently, in our N2pc analyses, target-
554    location congruent and incongruent distractor ERPs were averaged, as a function of the side
555    of distractor presentation.
556         *Electrical Neuroimaging of the N2pc component.* Our electrical neuroimaging
557    analyses separately tested response strength and topography in N2pc-like lateralised ERPs
558    (see e.g. Matusz et al., 2019b for a detailed, tutorial-like description of how electrical
559    neuroimaging measures can aid the study of attentional control processes). We assessed if
560    interactions between visual goals, multisensory salience and contextual factors 1)
561    modulated the distractor-elicited lateralised ERPs, and 2) if they do so by altering the
562    strength of responses within statistically indistinguishable brain networks and/or altering
563    the recruited brain networks.
564         *I. Lateralised analyses.* To test for the involvement of strength-based spatially-
565    selective mechanisms, we analysed Global Field Power (GFP) in lateralised ERPs. GFP is the
566    root mean square of potential [μV] across the entire electrode montage (see Lehmann &
567    Skrandies, 1980). To test for the involvement of network-related spatially-selective
568    mechanisms, we analysed stable patterns in ERP topography characterising different
569    experimental conditions using a clustering approach known as the Topographic Atomize and
570    Agglomerate Hierarchical Clustering (TAAHC). This topographic clustering procedure
571    generated sets of clusters of topographical maps that explained certain amounts of variance
572    within the group-averaged ERP data. Each cluster was labelled with a 'template map' that
573    represented the centroid of its cluster. The optimal number of clusters is one that explains
574    the largest global explained variance in the group-averaged ERP data with the smallest
575    number of template maps, and which we identified using the modified Krzanowski–Lai
576    criterion (Murray et al., 2008). In the next step, i.e., the so-called fitting procedure, the
577    single-subject data was 'fitted' back onto the topographic clustering results, such that each

578   datapoint of each subject's ERP data over a chosen time-window was labelled by the
579   template map with which it was best spatially correlated. This procedure resulted in a
580   number of timeframes that a given template map was present over a given time-window,
581   which durations (in milliseconds) we then submitted to statistical analyses described below.
582       In the present study, we conducted strength- and topographic analyses using the
583   same 4-way repeated-measures design as in the behavioural and canonical N2pc analyses,
584   on the lateralised whole-montage ERP data. Since the N2pc is a lateralised ERP, we first
585   conducted an electrical neuroimaging analysis of lateralised ERPs in order to uncover the
586   modulations of the N2pc by contextual factors. To obtain *global* electrical neuroimaging
587   measures of *lateralised* N2pc effects, we computed a difference ERP by subtracting the
588   voltages over the contralateral and ipsilateral hemiscalp, separately for each of the 4
589   distractor conditions. This resulted in a 59-channel difference ERP (as the midline electrodes
590   from the 129-electrode montage were not informative). Next, this difference ERP was
591   mirrored onto the other side of the scalp, recreating a "fake" 129 montage (with values on
592   midline electrodes now set to 0). It was on these mirrored "fake" 129-channel lateralised
593   difference ERPs that lateralised strength-based and topography-based electrical
594   neuroimaging analyses were performed. Here, GFP was extracted over the canonical 180–
595   300ms N2pc time-window and submitted to a 2 × 2 × 2 × 2 rmANOVA with factors Distractor
596   Colour (TCC vs. NCC), Distractor Modality (V vs. AV), as well as the two new factors,
597   Multisensory Relationship (Arbitrary vs. Congruent), and Distractor Onset (Distractor Onset;
598   Unpredictable vs. Predictable). Meanwhile, for topographic analyses, the "fake" 129-
599   channel data across the 4 Tasks were submitted to a topographic clustering over the entire
600   post-distractor period. Next, the data were fitted back over the 180-300ms period. Finally,
601   the resulting number of timeframes (in ms) was submitted to the same rmANOVA as the
602   GFP data above.
603       It remains unknown if the tested contextual factors modulate lateralised ERP
604   mechanisms at all. Given evidence that semantic information and temporal expectations
605   can modulate *nonlateralised* ERPs within the first 100–150ms post-stimulus (e.g., Dell'Acqua
606   et al., 2010; Dassanayake et al., 2016), we also investigated the influence of contextual
607   factors on nonlateralised voltage gradients, in an exploratory fashion. It must be noted that
608   ERPs are sensitive to the inherent physical differences in visual and audiovisual conditions.
609   Specifically, on audiovisual trials, the distractor-induced ERPs would be contaminated by
610   brain response modulations induced by sound processing, with these modulations visible in
611   our data already at 40ms post-distractor. Consequently, any direct comparison of visual-
612   only and audiovisual ERPs would index auditory processing per se and not capture of
613   attention by audiovisual stimuli. Such confounded sound-related activity is eliminated in the
614   canonical N2pc analyses through the contralateral-minus-ipsilateral subtraction. To
615   eliminate this confound in our electrical neuroimaging analyses here, we calculated
616   difference ERPs, first between TCCV and NCCV conditions, and then between TCCAV and
617   NCCAV conditions. Such difference ERPs, just as the canonical N2pc difference waveform,
618   subtract out the sound processing confound in visually-induced ERPs. As a result of those
619   difference ERPs, we removed factors Distractor Colour and Distractor Modality, and
620   produced a new factor, Target Difference (two levels: $D_{AV}$ [TCCAV − NCCAV difference] and
621   $D_V$ [TCCV − NCCV difference]), that indexed the enhancement of visual attentional control by
622   sound presence.
623       *II. Nonlateralised analyses.* All nonlateralised electrical neuroimaging analyses
624   involving context factors were based on the Target Difference ERPs. Strength-based

625     analyses, voltage and GFP data were submitted to 3-way rmANOVAs with factors:
626     Multisensory Relationship (Arbitrary vs. Congruent), Distractor Onset (Unpredictable vs.
627     Predictable), and Target Difference ($D_{AV}$ vs. $D_V$), and analysed using the STEN toolbox 1.0
628     (available for free at https://zenodo.org/record/1167723#.XS3lsi17E6h). Follow-up tests
629     involved further ANOVAs and pairwise *t*-tests. To correct for temporal and spatial
630     correlation (see Guthrie & Buchwald, 1991), we applied a temporal criterion of >15
631     contiguous timeframes, and a spatial criterion of >10% of the 129- channel electrode
632     montage at a given latency for the detection of statistically significant effects at an alpha
633     level of 0.05. As part of topography-based analyses, we segmented the ERP difference data
634     across the post-distractor and pre-target onset period (0 – 300ms from distractor onset). To
635     isolate the effects related to each of the two cognitive processes and reduce the complexity
636     of the performed analyses, we carried out two topographic clustering analyses. Topographic
637     clustering on nonlinear mechanisms contributing to TAC was based on the visual Target
638     Difference ERPs, while the clustering isolating MSE was based on difference ERPs resulting
639     from the subtraction of $D_{AV}$ and $D_V$. Thus, 4 group-averaged ERPs were submitted to both
640     clustering analyses, one for each of the context-related conditions. Next, the data were
641     fitted onto the canonical N2pc time-window (180–300ms) as well as other, earlier time-
642     periods, notably, also ones including time-periods highlighted by the GFP results as
643     representing significant condition differences. The resulting map presence (in ms) over the
644     given time-windows were submitted to 4-way rmANOVAs with factors: Multisensory
645     Relationship (Arbitrary vs. Congruent), Distractor Onset (Unpredictable vs. Predictable), and
646     Map (different numbers of maps, depending on the topographic clustering analyses and
647     time-windows within each clustering analyses), followed by post-hoc *t*-tests. Maps with
648     durations <15 contiguous timeframes were not included in the analyses. Unless otherwise
649     stated in the Results, map durations were statistically different from 0ms (as confirmed by
650     post-hoc one-sample t-tests), meaning that they were reliably present across the time-
651     windows of interest. Holm-Bonferroni corrections (Holm, 1979) were used to correct for
652     multiple comparisons between map durations. Comparisons passed the correction unless
653     otherwise stated.
654
655                                          **Results**
656
657     **Behavioural analyses**
658     ***Interaction of TAC and MSE with contextual factors***
659             To shed light on attentional control in naturalistic settings, we first tested whether
660     top-down visual control indexed by TAC interacted with contextual factors in behavioural
661     measures. First, our 2 × 2 × 2 × 2 rmANOVA confirmed the presence of TAC, via a main effect
662     of Distractor Colour, $F_{(1, 38)}$ = 340.4, $p < 0.001$, $\eta_p^2$ = 0.9, with TCC distractors (42ms), but not
663     NCC distractors (-1ms), eliciting reliable behavioural capture effects. Of central interest
664     here, the strength of TAC was dependent on whether the multisensory relationship within
665     the distractor involved mere simultaneity or semantic congruence. This was demonstrated
666     by a 2-way Distractor Colour × Multisensory Relationship interaction, $F_{(1, 38)}$ = 4.5, $p = 0.041$,
667     $\eta_p^2$ = 0.1 (Figure 2). This effect was driven by behavioural capture effects elicited by TCC
668     distractors being reliably larger for the Arbitrary (45ms) than for the Congruent (40ms)
669     condition, $t_{(38)}$ = 1.9, $p = 0.027$. NCC distractors showed no evidence of Multisensory
670     Relationship modulation (Arbitrary vs. Congruent, $t_{(38)}$ = 1, $p = 0.43$). Contrastingly, TAC
671     showed no evidence of modulation by predictability of the distractor onset (no 2-way

672    Distractor Colour × Distractor Onset interaction, $F_{(1, 38)}$ = 2, $p$ = 0.16). Thus, visual feature-
673    based attentional control interacted with the contextual factor of distractor semantic
674    congruence, but not distractor temporal predictability.
675        Next, we investigated potential interactions of multisensory enhancements with
676    contextual factors. Expectedly, there was behavioural MSE (a significant main effect of
677    Distractor Modality, $F_{(1, 38)}$=13.5, $p$=0.001, $\eta_p^2$=0.3), where visually-elicited behavioural
678    capture effects (18ms) were enhanced on AV trials (23ms). Unlike TAC, this MSE effect
679    showed no evidence of interaction with either of the two contextual factors (Distractor
680    Modality x Multisensory Relationship interaction, $F$<1; Distractor Modality x Distractor
681    Onset interaction: $n.s.$ trend, $F_{(1, 38)}$=3.6, $p$=0.07, $\eta_p^2$= 0.1). Thus, behaviourally, Multisensory
682    enhancement of attentional capture was not modulated by distractors' semantic
683    relationship nor its temporal predictability. We have also observed other, unexpected
684    effects, but as these were outside of the focus of the current paper, which aims to elucidate
685    the interactions between visual (goal-based) and multisensory (salience-driven) attentional
686    control and contextual mechanisms, we describe them only in SOMs.
687    
688    
689                                    ** FIGURE 2 HERE **
690    
691    
692    
693    **ERP analyses**
694    ***Lateralised (N2pc-like) brain mechanisms***
695    We next investigated the type of brain mechanisms that underlie interactions between
696    more traditional attentional control (TAC, MSE) and contextual control over attentional
697    selection. Our analyses on the lateralised responses, spanning both a canonical and EN
698    framework, revealed little evidence for a role of spatially-selective mechanisms in
699    supporting the above interactions. Both canonical N2pc and electrical neuroimaging
700    analyses confirmed the presence of TAC (see Fig. 3 for N2pc waveforms across the four
701    distractor types). However, TAC did not interact with either of the two contextual factors.
702    Lateralised ERPs also showed no evidence for sensitivity to MSE nor for interactions
703    between MSE and any contextual factors. Not even the main effects of Multisensory
704    Relationship and Distractor Onset[6] were present in lateralised responses (See SOMs for full
705    description of the results of lateralised ERP analyses).
706    
707    
708                                    ** FIGURE 3 HERE **
709    
710    
711    ***Nonlateralised brain mechanisms***
712    A major part of our analyses focused on understanding the role of nonlateralised ERP
713    mechanisms in the interactions between visual goals (TAC), multisensory salience (MSE) and
714    contextual control. To remind the reader, to prevent nonlateralised ERPs from being

---

[6] Any ERP results related to Distractor Onset are unlikely to be confounded by shifted baseline due to potential dominance of one ISI type (100ms, 250ms, 450ms) over others, as no such dominance was identified in a subsample of data.

16

715　　confounded by the presence of sound on AV trials, we based our analyses here on the
716　　difference ERPs indexing visual attentional control under sound absence vs. presence. That
717　　is, we calculated ERPs of the difference between TCCV and NCCV conditions, and between
718　　TCCAV and NCCAV conditions ($D_V$ and $D_{AV}$ levels, respectively, of the Target Difference
719　　factor). We focus the description of these results on the effects of interest (see SOMs for full
720　　description of results).
721　　　　　The 2 × 2 × 2 (Multisensory Relationship × Distractor Onset × Target Difference)
722　　rmANOVA on electrode-wise voltage analyses revealed a main effect of Target Difference at
723　　53–99ms and 141–179ms, thus both at early, perception-related, and later, attentional
724　　selection-related latencies (reflected by the N2pc). Across both time-windows, amplitudes
725　　were larger for $D_{AV}$ (TCCAV − NCCAV difference) than for $D_V$ (TCCV − NCCV difference). This
726　　effect was further modulated, evidenced by a 2-way Target Difference × Multisensory
727　　Relationship interaction, at the following time-windows: 65–103ms, 143–171ms, and 194–
728　　221ms (all $p$'s < 0.05). The interaction was driven by Congruent distractors showing larger
729　　amplitudes for $D_{AV}$ than $D_V$ within all 3 time-windows (65–97ms, 143–171ms, and 194–
730　　221ms; all $p$'s < 0.05). No similar differences were found for Arbitrary distractors, and there
731　　were no other interactions that passed the temporal and spatial criteria for multiple
732　　comparisons of >15 contiguous timeframes and >10% of the 129- channel electrode
733　　montage.
734
735　　***Interaction of TAC with contextual factors.*** We next used electrical neuroimaging analyses
736　　to investigate the contribution of the strength- and topography-based nonlateralised
737　　mechanisms to the interactions between TAC and contextual factors.
738　　　　　*Strength-based brain mechanisms.* A 2 × 2 × 2 Target Difference × Multisensory
739　　Relationship × Distractor Onset rmANOVA on the GFP mirrored the results of the electrode-
740　　wise analysis on ERP voltages by showing a main effect of Target Difference spanning a large
741　　part of the first 300ms post-distractor both before and in N2pc-like time-windows (19–
742　　213ms, 221–255ms, and 275–290ms). Like in the voltage waveform analysis, the GFP was
743　　larger for $D_{AV}$ than $D_V$ (all $p$'s < 0.05). In GFP, Target Difference interacted both with
744　　Multisensory Relationship (23–255ms) and separately with Distractor Onset (88–127ms; see
745　　SOMs for full description). Notably, there was a 3-way Target Difference × Multisensory
746　　Relationship × Distractor Onset interaction, spanning 102–124ms and 234–249ms. We
747　　followed up this interaction with a series of post-hoc tests to gauge the modulations of TAC
748　　(and MSE, see below) by the two contextual factors.
749　　　　　In GFP, Multisensory Relationship and Distractor Onset interacted independently of
750　　Target Difference in the second time-window, which results we describe in SOMs. To gauge
751　　differences in the strength of TAC in GFP across the 4 contexts (i.e., Arbitrary Unpredictable,
752　　Arbitrary Predictable, Congruent Unpredictable, and Congruent Predictable), we focused
753　　the comparisons on only visually-elicited target differences (to minimise any potential
754　　confounding influences from sound processing) across the respective levels of the 2
755　　contextual factors. The weakest GFPs were observed for Arbitrary Predictable distractors
756　　(Figure 4A). They were weaker than GFPs elicited for Arbitrary Unpredictable distractors
757　　(102–124ms and 234–249ms), and Predictable Congruent distractors (only in the later
758　　window, 234–249ms).
759　　　　　*Topography-based brain mechanisms.* We focused the topographic clustering of the
760　　TAC-related topographic activity on the whole 0–300ms post-distractor time-window
761　　(before the target onset), which revealed 10 clusters that explained 82% of the global

762    explained variance within the visual-only ERPs. This time-window of 29–126ms post-
763    distractor was selected on based on the GFP peaks, which are known to correlate with
764    topographic stability (Lehmann 1987; Brunet et al. 2011), and in some conditions, based on
765    the fact that specific template was dominated responses in group-averaged data from given
766    conditions, e.g., Arbitrary Unpredictable and Congruent Unpredictable conditions, but not
767    for other conditions. This was confirmed by our statistical analyses, with a 2 × 2 × 5
768    rmANOVA over the 29–126ms post-distractor time-window, which revealed a 3-way
769    Multisensory Relationship × Distractor Onset × Map interaction, $F_{(3.2,122)}$ = 5.3, $p$ = 0.002, $\eta_p^2$
770    = 0.1.

771          Follow-up tests in the 29–126ms time-window focused on maps differentiating
772    between the 4 contexts as a function of the two contextual factors (results of follow-up
773    analyses as a function of Multisensory Relationship and Distractor Onset are visible in Figure
774    4B in leftward panel and rightward panel, respectively). These results confirmed that
775    context altered the processing of distractors from early on. The results also confirmed the
776    clustering that the context did so by engaging different networks for most of the different
777    combinations of Multisensory Relationship and Distractor Onset: Arbitrary Unpredictable -
778    Map A2, Congruent Unpredictable - Map A5, as well as for Arbitrary Predictable - Map A1
779    (no map predominantly involved in the responses for Congruent Predictable).

780          *Arbitrary Predictable distractors,* which elicited the weakest GFP, recruited
781    predominantly Map A1 (37ms) during processing. This map was more involved in the
782    processing of those distractors vs. Congruent Predictable distractors (21ms), $t_{(38)}$ = 2.7, $p$ =
783    0.013 (Fig.4B bottom panel).

784          *Arbitrary Unpredictable* distractors largely recruited Map A2 (35ms) during
785    processing. This map was more involved in the processing of these distractors vs. Arbitrary
786    Predictable distractors (18ms), $t_{(38)}$ = 2.64, $p$ = 0.012 (Fig.4B top leftward panel), as well as
787    Congruent Unpredictable distractors (14ms), $t_{(38)}$ = 3.61, $p$ < 0.001 (Fig.4B top rightward
788    panel).

789          *Congruent Unpredictable* distractors principally recruited Map A5 (34ms) during
790    processing, which was more involved in the processing of these distractors vs. Congruent
791    Predictable distractors (19ms) distractors, $t_{(38)}$= 2.7, $p$ = 0.039 (Fig.4B middle leftward
792    panel), as well as Arbitrary Unpredictable (12ms) distractors, $t_{(38)}$ = 3.7, $p$ <0.001 (Fig.4B
793    middle rightward panel).

794          *Congruent Predictable* distractors recruited different template maps during
795    processing, where Map A2 was more involved in responses to those distractors (25ms) vs.
796    Congruent Unpredictable distractors (14ms), $t_{(38)}$ = 2.17, $p$ = 0.037, but not other distractors,
797    $p$'s>0.2 (Fig.4B top leftward panel).

798

799    **Interaction of MSE with contextual factors.** We next analysed the strength- and
800    topography-based nonlateralised mechanisms contributing to the interactions between
801    MSE and contextual factors.

802          *Strength-based brain mechanisms.* To gauge the AV-induced enhancements between
803    $D_{AV}$ and $D_V$ across the 4 contexts, we explored the abovementioned 2 × 2 × 2 GFP interaction
804    using a series of simple follow-up post-hoc tests. We first tested if response strength
805    between $D_{AV}$ and $D_V$ was reliably different within each of the 4 contextual conditions. AV-
806    induced ERP responses were enhanced (i.e., larger GFP for $D_{AV}$ than $D_V$ distractors) for both
807    Predictable and Unpredictable Congruent distractors, across both earlier and later time-
808    windows. Likewise, AV enhancements were also found for Arbitrary Predictable distractors,

18

809  but only in the earlier (102–124ms) time-window. Unpredictable distractors showed similar
810  GFP across $D_{AV}$ and $D_V$ trials. Next, we compared the AV-induced MS enhancements across
811  the 4 contexts, by creating ($D_{AV}$ minus $D_V$) difference ERPs or each context. AV-induced
812  enhancements were weaker for Predictable Arbitrary distractors than Predictable
813  Congruent distractors (102–124ms and 234–249ms; Figure 5A).
814
815
816                              ** FIGURE 4 HERE **
817
818
819          *Topography-based brain* mechanisms. We then used the difference ($D_{AV}$ minus $D_V$)
820  difference ERPs (as in the second part of the GFP analyses) to focus the topographic
821  clustering selectively on the MSE-related topographic activity. This clustering, carried out on
822  the 0–300ms post-distractor and pre-target time-window, revealed 7 clusters that explained
823  78% of the global explained variance within the AV-V target difference ERPs.
824          In this topographic clustering there were multiple GPF peaks, with elongated near-
825  synchronous periods of time where different maps were suggested to be present across the
826  four distractor conditions in the group-averaged data. One of those maps (Map B3) was first
827  present in the two congruent distractor conditions, to then become absent and reappear
828  again. In the view of this patterning, we decided to fit the group-average data from these
829  three subsequent time-windows to single-subject data: 35–110ms, 110– 190ms, and 190–
830  300ms. To foreshadow the results, in the first and third time-windows the MSE-related
831  template maps were modulated only by Multisensory Relationship, while in the middle
832  time-window – by both Multisensory Relationship and Distractor Onset.
833          In the first, 35–110ms time-window, the modulation of map presence by
834  Multisensory Relationship was evidenced by a 2-way Map × Multisensory Relationship
835  interaction, $F_{(2.1,77.9)}$ = 9.2, $p$ < 0.001, $\eta_p^2$ = 0.2. This effect was driven by one map (map B3)
836  that, in this time-window, predominated responses to Congruent (42ms) vs. Arbitrary
837  (25ms) distractors, $t_{(38)}$ = 4.3, $p$ = 0.02, whereas another map (map B5) dominated responses
838  to Arbitrary (33ms) vs. Congruent (18ms) distractors, $t_{(38)}$ = 4, $p$ = 0.01 (Figure 5B top and
839  upper leftward panels, respectively).
840          In the second, 110–190ms time-window, map presence was modulated by both
841  contextual factors, with a 3-way Map × Multisensory Relationship × Distractor Onset
842  interaction, $F_{(2.6,99.9)}$ = 3.7, $p$ = 0.02, $\eta_p^2$ = 0.1 (just as it did for TAC). We focused follow-up
843  tests in that time-window again on maps differentiating between the 4 conditions, as we did
844  for the 3-way interaction for TAC (results of follow-ups as a function of Multisensory
845  Relationship and Distractor Onset are visible in Figure 5B, middle upper and lower panels,
846  respectively). Context processes again interacted to modulate the processing of distractors,
847  although now they did so after the first 100ms. They did so again by engaging different
848  networks for different combinations of Multisensory Relationship and Distractor Onset:
849  Arbitrary Predictable distractors - Map B1, Arbitrary Unpredictable distractors - Map B5,
850  Congruent Unpredictable distractors - Map B6, and now also Congruent Predictable
851  distractors - Map B3.
852
853
854                              ** FIGURE 5 HERE **
855

19

856
857  *Arbitrary Predictable* distractors, which again elicited the weakest GFP, during
858  processing mainly recruited map B1 (35ms). This map dominated responses to these
859  distractors vs. Arbitrary Unpredictable distractors (18ms, $t_{(38)}$ = 2.8, $p$ = 0.01; Figure 5B
860  upper panel), as well as Congruent Predictable distractors (17ms, $t_{(38)}$ = 2.8, $p$ = 0.006; Figure
861  5B lower panel).
862  *Arbitrary Unpredictable* distractors largely recruited during processing one map,
863  Map B5 (33ms). Map B5 was more involved in responses to these distractors vs. Arbitrary
864  Predictable distractors (17ms, $t_{(38)}$ = 2.6, $p$ = 0.042; Figure 5B  upper panel), as well as vs.
865  Congruent Unpredictable distractors (13ms, $t_{(38)}$ = 3.4, $p$ = 0.002; Figure 5B bottom panel).
866  *Congruent Unpredictable* distractors principally recruited during processing Map B6
867  (37ms). Map B6 was more involved in responses to these distractors vs. Congruent
868  Predictable distractors (21ms, $t_{(38)}$ = 2.5, $p$ = 0.02), and vs. Arbitrary Unpredictable
869  distractors (24ms, $t_{(38)}$ = 2.3, $p$ = 0.044).
870  *Congruent Predictable* distractors mostly recruited during processing Map B3 (25ms).
871  Map B3 was more involved in responses to these distractors vs. Predictable Arbitrary
872  distractors (8ms, $t_{(38)}$ = 2.2, $p$ = 0.005), and, at statical-significance threshold level, vs.
873  Congruent Unpredictable distractors (12ms, $t_{(38)}$ = 2.2, $p$ = 0.0502).
874  In the third, 190–300ms time-window, the 2-way Map × Multisensory Relationship
875  interaction was reliable at $F_{(3.2, 121.6)}$ = 3.7, $p$ = 0.01, $\eta_p^2$ = 0.1. Notably, the same map as
876  before (map B3) was more involved, at a non-statistical trend level, in the responses to
877  Congruent (50ms) vs. Arbitrary distractors (33ms), $t_{(38)}$ = 3.6, $p$ = 0.08, and another map
878  (map B1) predominated responses to Arbitrary (25ms) vs. Congruent (14ms) distractors, $t_{(38)}$
879  = 2.3, $p$ = 0.02 (Figure 5B rightward panel).
880
881  **Discussion**
882
883  Attentional control is necessary to cope with the multitude of stimulation in everyday
884  situations. However, in such situations, the observer's goals and stimuli's salience routinely
885  interact with contextual processes, yet such multi-pronged interactions between control
886  processes have never been studied. Below, we discuss our findings on how visual and
887  multisensory attentional control interact with distractor temporal predictability and
888  semantic relationship. We then discuss the spatiotemporal dynamics in nonlateralised brain
889  mechanisms underlying these interactions. Finally, we discuss how our results enrich the
890  understanding of attentional control in real-world settings.
891
892  **Interaction of task-set contingent attentional capture with contextual control**
893  Visual control interacted most robustly with stimuli's semantic relationship. Behaviourally,
894  *target-matching* visual distractors captured attention more strongly when they were
895  arbitrarily connected than semantically congruent. This was accompanied by a cascade of
896  modulations of nonlateralised brain responses, spanning both the attentional selection,
897  N2pc-like stage and much earlier, perceptual stages. Arbitrary distractors, but only
898  predictable ones, first recruited one particular brain network (Map A1), to a larger extent
899  than predictable semantically congruent distractors, and did so early on (29–126ms post-
900  distractor). Arbitrary predictable distractors elicited also suppressed responses, in the later
901  part of this early time-window (102–124ms; where they elicited the weakest responses). In
902  the later, N2pc-like (234–249ms) time-window, responses to arbitrary predictable

20

903   distractors were again weaker, now compared to semantically congruent predictable
904   distractors.
905         This cascade of network- and strength-based modulations of nonlateralised brain
906   responses might epitomise a potential brain mechanism for interactions between visual top-
907   down control and multiple sources of contextual control, as they are consistent with existing
908   literature. The discovered early (~30-100ms) topographic modulations for predictable target-
909   matching (compared to unpredictable) distractors is consistent with predictions attenuating
910   the earliest visual perceptual stages (C1 component, ~50–100ms post-stimulus;
911   Dassanayake et al. 2016). The subsequent, mid-latency response suppressions (102–124ms,
912   where we found also topographic modulations) for predictable distractors are in line with N1
913   attenuations for self-generated sounds (Baess et al. 2011; Klaffehn et al. 2019), and the
914   latencies where the brain might promote the processing of unexpected events (Press et al.
915   2020). Notably, these latencies are also in line with the onset (~115ms post-stimulus) of the
916   goal-based suppression of salient visual distractors (here: presented simultaneously with
917   targets), i.e., distractor positivity (Pd; Sawaki & Luck 2010). Finally, the response
918   suppressions we found at later, N2pc-like, attentional selection stages (234–249ms), are
919   also consistent with some extant (albeit scarce) literature. Van Moorselaar and Slagter
920   (2019) showed that when such salient visual distractors appear in predictable locations,
921   they elicit the N2pc but no longer a (subsequent, post-target) Pd, suggesting that once the
922   brain learns the distractor's location, it can suppress it without the need for active
923   inhibition. More recently, van Moorselaar et al. (2020b) showed that the representation of
924   the predictable distractor feature could be decoded already from pre-stimulus activity.
925   While our paradigm was not optimised for revealing such effects, pre-stimulus mechanisms
926   could indeed explain our early-onset (~30ms) context-elicited neural effects. The robust
927   response suppressions for predictable stimuli are also consistent with recent proposals for
928   interactions between predictions and auditory attention. Schröger et al. (2015) suggested
929   that greater attention is deployed to more "salient" stimuli, i.e., those for which a prediction
930   is missing, so that the predictive model can be reconfigured to encompass such predictions
931   in the future. This reconfiguration, in turn, requires top-down goal-based attentional
932   control. Our results extend this model to the visual domain. Our findings involving the
933   response modulation cascade and behavioural benefits may also support the Schröger et
934   al.'s tenet that different, but connected, predictive models exist at different levels of the
935   cortical hierarchy.
936         These existing findings jointly strengthen our interpretations that goal-based top-
937   down control utilises contextual information to alter visual processing from very early on in
938   life. Our findings also extend the extant ideas in several ways. First, they show that in
939   context-rich settings (i.e., involving multiple sources of contextual control), goal-based
940   control will use both stimulus-related predictions and stimulus meaning to facilitate task-
941   relevant processing. Second, context information modulates not only early, pre-stimulus
942   and late, attentional stages, but also early *stimulus-elicited* responses. Third, our findings
943   also suggest candidate mechanisms for supporting interactions between goal-based control
944   and multiple sources of contextual information. Namely, context will modulate the early
945   stimulus processing by recruiting distinct brain networks for stimuli representing different
946   contexts, e.g., the brain networks recruited by predictable distractors differed for arbitrarily
947   linked and semantically congruent stimuli (Map A1 and A2, respectively). Also, the distinct
948   network recruitment might lead to the suppressed (potentially more efficient; c.f. repetition
949   suppression, Grill-Spector et al. 2006) brain responses. These early response attenuations

21

950 will extend also to later stages, associated with attentional selection. Thus, it is the early
951 differential brain network recruitment that might trigger a cascade of spatiotemporal brain
952 dynamics leading effectively to the stronger behavioural capture, here for predictable
953 (arbitrary) distractors. However, for distractors, these behavioural benefits may be most
954 robust for arbitrary target-matching stimuli (as opposed to semantically congruent), with
955 prediction-based effects are less apparent.
956
957 **Interaction of multisensory enhancement of attentional capture with contextual control**
958 Across brain responses, multisensory-induced processes interacted with both contextual
959 processes. To measure effects related to multisensory-elicited modulations and to its
960 interactions with contextual information, we analysed AV–V differences within the Target
961 Difference ERPs.
962 The interactions between multisensory modulations and context processes were
963 also instantiated via an early-onset cascade of strength- and **topographic** (network-based)
964 nonlateralised brain mechanisms. This cascade again started early (now 35–110ms post-
965 distractor). A separate topographic clustering analysis revealed that in the multisensory-
966 modulated responses the brain first distinguished only between semantically congruent and
967 arbitrarily linked distractors. These distractors recruited predominantly different brain
968 networks (Map B3 and B5, respectively). Around the end of these topographic, network-
969 based modulations, at 102–124ms, multisensory-elicited brain responses were also
970 modulated in their strength. Arbitrary predictable distractors again triggered weaker
971 responses, now compared to semantically congruent predictable distractors. Multisensory-
972 elicited responses predominantly recruited distinct brain networks for the four distractor
973 types from 110ms until 190ms post-distractor, thus spanning stages linked to perception
974 and attentional selection. Here, maps B3 and B5 were now recruited for responses to
975 semantically congruent predictable and arbitrary unpredictable distractors, respectively.
976 Meanwhile, maps B1 and B6 were recruited for arbitrary predictable and semantically
977 congruent unpredictable distractors, respectively. In the subsequent time-window (190–
978 300ms) that mirrors the time-window used in the canonical N2pc analyses, multisensory-
979 related responses again recruited different brain networks. There, Map B3 (previously:
980 Congruent Predictable distractors) again was predominantly recruited by semantically
981 congruent over arbitrary distractors, and now Map B1 (previously: Arbitrary Predictable
982 distractors) - for arbitrary distractors over congruent ones. In the middle of this time-
983 window (234–249ms), responses differed in their strength, with predictable arbitrary
984 distractors eliciting weaker responses compared to semantically congruent predictable
985 distractors.
986 To summarise, distractors' semantic relationship played a dominant (but not
987 absolute) role in interactions between multisensory-elicited and contextual processes. The
988 AV–V difference ERPs were modulated exclusively by multisensory relationships both in the
989 earliest, perceptual (35–110ms) time-window and latest, N2pc-like (190–300ms) time-
990 window linked to attentional selection. At both stages, distinct brain networks were
991 recruited predominantly by semantically congruent and arbitrary distractors. These results
992 suggest that from early perceptual stages the brain "relays" the processing of (multisensory)
993 stimuli as a function of them containing meaning (vs. lack thereof) for the observer up to
994 stages of attentional selection. Notably, the same brain network (Map B3) supported
995 multisensory processing of semantically congruent distractors across both time-windows,
996 while different networks were recruited by arbitrarily linked distractors.

22

997        Thus, a single network might be recruited for processing meaningful multisensory
998 stimuli. In light of our behavioural results, this brain network could be involved in
999 suppressing behavioural attentional capture for semantically congruent (over arbitrarily
1000 linked) distractors by top-down goal-driven attentional control. This idea is supported by the
1001 interactions between distractors' multisensory-driven modulations, their multisensory
1002 relationship, and their temporal predictability in the second, 110–190ms time-window.
1003 Therein, the same "semantic" Map B3 was still present, albeit now recruited for responses
1004 to semantically congruent (over arbitrary) *predictable* distractors. Based on existing
1005 evidence that predictions are used in service of goal-based behaviour (Schröger et al. 2015;
1006 van Moorselaar et al. 2020a; Matusz et al. 2016), one could argue that the brain network
1007 reflected by Map B3 might play a role in integrating contextual information across both
1008 predictions and meaning (though mostly meaning, as it remained recruited by semantically
1009 congruent distractors throughout the distractor-elicited response). The activity of this
1010 network might have contributed to the overall stronger brain responses (indicated by GFP
1011 results) to semantically congruent multisensory stimuli, which in turn contributed to the null
1012 behavioural multisensory enhancements of behavioural indices of attentional capture.
1013 While these are the first results of this kind, they open an exciting possibility that surface-
1014 level EEG/ERP studies can reveal the network- and strength-related brain mechanisms
1015 (potentially a single network for "gain control" up-modulation) by which goal-based
1016 processes control (i.e., suppress) multisensorily-driven enhancements of attentional
1017 capture.
1018
1019 **Towards understanding how we pay attention in naturalistic settings**
1020 It is now relatively well-established that the brain facilitates goal-directed processing (from
1021 perception to attentional selection) via processes based on observer's goals (e.g. Folk et al.
1022 1992; Desimone & Duncan 1995), predictions about the outside world (Summerfield &
1023 Egner 2009; Schröger et al. 2015; Press et al. 2020), and long-term memory contents
1024 (Summerfield et al. 2006; Peelen & Kastner 2014). Also, multisensory processes are
1025 increasingly recognised as an important source of bottom-up, attentional control (e.g.
1026 Spence & Santangelo 2007; Matusz & Eimer 2011; Matusz et al. 2019a; Fleming et al. 2020).
1027 By studying these processes largely in isolation, researchers clarified how they support goal-
1028 directed behaviour. However, in the real world, observer's goals interact with multisensory
1029 processes and multiple types of contextual information. Our study sheds first light on this
1030 "naturalistic attentional control".
1031        Understanding of attentional control in the real world has been advanced by
1032 research on feature-related mechanisms (Theeuwes 1991; Folk et al. 1992; Desimone &
1033 Duncan 1995; Luck et al. 2020), which support attentional control where target location
1034 information is missing. Here, we aimed to increase the ecological validity of this research by
1035 investigating how visual feature-based attention (as indexed by TAC) transpires in context-
1036 rich, multisensory settings (see SOMs for a discussion of our replication of TAC). Our findings
1037 of reduced capture for semantically congruent than artificially linked target-colour matching
1038 distractors is novel and important, as they suggest stimuli's meaning is also utilised to
1039 suppress attention (to distractors). Until now, known benefits of meaning were limited to
1040 target selection (Thorpe et al. 1996; Iordanescu et al. 2008; Matusz et al. 2019a). Folk et al.
1041 (1992) famously demonstrated that attentional capture by distractors is sensitive to the
1042 observer's goals; we reveal that distractor's meaning may serve as a second source of goal-
1043 based attentional control. This provides a richer explanation for how we stay focused on

1044    task in everyday situations, despite many objects matching attributes of our current
1045    behavioural goals.
1046        To summarise, in the real world, attention should be captured more strongly by
1047    stimuli that are unpredictable (Schröger et al. 2015), but also by those unknown or without
1048    a clear meaning. On the other hand, stimuli with high strong spatial and/or temporal
1049    alignment across senses (and so stronger bottom-up salience) may be more resistant to
1050    such goal-based attentional control (suppression), as we have shown here (multisensory
1051    enhancement of attentional capture; see also Santangelo & Spence 2007; Matusz & Eimer
1052    2011; van der Burg et al. 2011; Turoman et al. 2021a; Fleming et al. 2020). As multisensory
1053    distractors captured attention more strongly even in current, context-rich settings, this
1054    confirms the importance of multisensory salience as a source of *potential* bottom-up
1055    attentional control in naturalistic environments (see SOMs for a short discussion of this
1056    replication).
1057        The investigation of brain mechanisms underlying known EEG/ERP correlates (N2pc,
1058    for TAC) via advanced multivariate analyses has enabled us to provide a comprehensive,
1059    novel account of attentional control in a multi-sensory, context-rich setting. Our results
1060    jointly support the primacy of goal-based control in naturalistic settings. Multisensory
1061    semantic congruence reduced behavioural attentional capture by target-matching colour
1062    distractors compared to arbitrarily linked distractors. Context modulated nonlateralised
1063    brain responses to target-related (TAC) distractors via a cascade of strength- and topographic
1064    mechanisms from early (~30ms post-distractor) to later, attentional selection stages. While
1065    these results are first of this kind and need replication, they suggest that context-based
1066    goal-directed modulations of distractor processing "snowball" from early stages (potentially
1067    involving pre-stimulus processes, e.g. van Moorselaar & Slagter, 2020) to control
1068    behavioural attentional selection. Responses to predictable arbitrary (target-matching)
1069    distractors revealed by our electrical neuroimaging analyses might have driven the larger
1070    behavioural capture for arbitrary than semantically congruent distractors. The former
1071    engaged distinct brain networks and triggered the weakest and potentially most efficient
1072    (Grill-Spector et al. 2006) responses. One reason for the absence of such effects in
1073    behavioural measures is the small magnitude of behavioural effects: while the TAC effect is
1074    ~50ms, both MSE effect and semantically-driven suppression were small, at around ~5ms.
1075    This may also be the reason why context-driven effects were absent in behavioural
1076    measures of multisensory enhancement of attentional capture, despite involving a complex,
1077    early-onsetting cascade of strength- and topographic modulations.
1078        Our results point to a potential brain mechanism by which semantic relationships
1079    influence goal-directed behaviour towards task-irrelevant information. Namely, our
1080    electrical neuroimaging analyses of surface-level EEG identified a brain network that is
1081    recruited by semantically congruent stimuli at early, perceptual stages, and that remains
1082    active at N2pc-like, attentional selection stages. While remaining cautious when interpreting
1083    our results, this network might have contributed to the consistently enhanced AV-induced
1084    responses for semantically congruent multisensory distractors. These enhanced brain
1085    responses together with the concomitant *suppressed behavioural attention* effects are
1086    consistent with a "gain control" mechanism, in the context of distractor processing (e.g.
1087    Sawaki & Luck 2010; Luck et al. 2020). Our results reveal that such "gain control", at least in
1088    some cases, operates by relaying processing of certain stimuli to distinct brain networks. We
1089    have purported the existence of such a "gain control" mechanism in a different study on
1090    (top-down) multisensory attention (e.g. Matusz et al. 2019c). While these are merely

24

1091   speculations that would require source estimations to be supported, the enhanced
1092   responses to meaningful distractors may thus reflect enhanced goal-based control over
1093   those stimuli. Such a process could potentially recruit a network involving the anterior
1094   hippocampus and putamen, which help maintain active representations of task-relevant
1095   information while updating the representation of to-be-suppressed information (McNab &
1096   Klingberg 2008; Sadeh et al. 2010; Jiang et al. 2015). Our electrical neuroimaging analyses of
1097   the surface-level N2pc data (see also Matusz et al. 2019c; Turoman et al. 2021a) might have
1098   potentially revealed when and how such memory-related brain networks modulate
1099   attentional control over task-irrelevant stimuli.
1100
1101   **N2pc as an index of attentional control**
1102   We have previously discussed the limitations of canonical N2pc analyses in capturing
1103   neurocognitive mechanisms by which visual top-down goals and multisensory bottom-up
1104   salience simultaneously control attention selection (Matusz et al. 2019b). The mean N2pc
1105   amplitude modulations are commonly interpreted as "gain control", but they can be driven
1106   by both strength- (i.e., "gain") and topographic (network-based) mechanisms. Canonical N2pc
1107   analyses cannot distinguish between those two brain mechanisms. Contrastingly, Matusz et
1108   al. (2019b) have shown evidence for both brain mechanisms underlying N2pc-like
1109   responses. These and other results of ours (Turoman et al. 2021a, 2021b) provided evidence
1110   from surface-level data for different brain sources contributing to the N2pc's, a finding that
1111   has been previously shown only in *source*-level data (Hopf et al. 2000). These findings point
1112   to a certain limitation of the N2pc (canonically analysed), which is an EEG *correlate* of
1113   attentional selection, but where other analytical approaches are necessary to reveal brain
1114   mechanisms of attentional selection.
1115        Here, we have shown that the lateralised, spatially-selective brain mechanisms,
1116   approximated by the N2pc and revealed by electrical neuroimaging analyses are limited in
1117   how they contribute to attentional control in some settings. Rich, multisensory, and
1118   context-laden influences over goal-based top-down attention are, in our current paradigm,
1119   not captured by such lateralised mechanisms. In contrast, nonlateralised (or at least
1120   *relatively less* lateralised, see Figures 4 and 5) brain networks seem to support such
1121   interactions for visual and multisensory distractors - from early on, leading to attentional
1122   selection. We nevertheless want to reiterate that paradigms that can gauge N2pc offer an
1123   important starting point for studying attentional control in less traditional multisensory
1124   and/or context-rich settings. There, multivariate analyses, and an electrical neuroimaging
1125   framework in particular, might be useful in readily revealing new mechanistic insights into
1126   attentional control.
1127
1128   **Broader implications**
1129   Our findings are important to consider when aiming to study attentional control, and
1130   information processing more generally, in naturalistic settings (e.g., while viewing movies,
1131   listening to audiostories) and veridical real-world environments (e.g. the classroom or the
1132   museum). Additionally, conceptualisations of ecological validity (Peelen et al. 2014; Shamay-
1133   Tsoory & Mendelsohn 2019; Vanderwal et al. 2019; Eickhoff et al. 2020; Cantlon 2020)
1134   should go beyond traditionally invoked components (e.g., observer's goals, context,
1135   socialness) to encompass contribution of multisensory processes. For example, naturalistic
1136   studies should compare unisensory and multisensory stimulus/material formats, to
1137   measure/estimate the contribution of multisensory-driven bottom-up salience to the

25

1138    processes of interest. More generally, our results highlight that hypotheses about how
1139    neurocognitive functions operate in everyday situations can be built already in the
1140    laboratory, if one manipulates systematically, together and across the senses, goals,
1141    salience, and context (van Atteveldt et al. 2018; Matusz et al. 2019c). Such a cyclical
1142    approach (Matusz et al. 2019a; see also Naumann et al. 2020 for a new tool to measure
1143    ecological validity of a study) involving testing of hypotheses across laboratory and veridical
1144    real-world settings could be highly promising for successfully bridging the two, typically
1145    separately pursued types of research. As a result, such an approach could create more
1146    complete theories of naturalistic attentional control.

1147                                    **References**
1148
1149   Alais, D., Newell, F., & Mamassian, P. (2010). Multisensory processing in review: from
1150         physiology to behaviour. *Seeing and perceiving*, *23*(1), 3-38.

1151   Baess, P., Horváth, J., Jacobsen, T., & Schröger, E. (2011). Selective suppression of
1152         self-initiated sounds in an auditory stream: An ERP study. *Psychophysiology*, *48*(9),
1153         1276-1283.

1154   Bevilacqua, D., Davidesco, I., Wan, L., Chaloner, K., Rowland, J., Ding, M., ... & Dikker, S.
1155         (2019). Brain-to-brain synchrony and learning outcomes vary by student–teacher
1156         dynamics: Evidence from a real-world classroom electroencephalography
1157         study. *Journal of cognitive neuroscience*, *31*(3), 401-411.

1158   Biasiucci, A., Franceschiello, B., & Murray, M. M. (2019). Electroencephalography. *Current,*
1159         *29(3)*, R80-R85.

1160   Brunet, D., Murray, M. M., & Michel, C. M. (2011). Spatiotemporal analysis of multichannel
1161         EEG: CARTOOL. *Computational intelligence and neuroscience*, *2011*.

1162   Burra, N., & Kerzel, D. (2013). Attentional capture during visual search is attenuated by
1163         target predictability: Evidence from the N2pc, Pd, and topographic
1164         segmentation. *Psychophysiology*, *50*(5), 422-430.

1165   Cantlon, J. F. (2020). The balance of rigor and reality in developmental
1166         neuroscience. *NeuroImage*, *216*, 116464.

1167   Cappe, C., Thut, G., Romei, V., & Murray, M. M. (2010). Auditory–visual multisensory
1168         interactions in humans: timing, topography, directionality, and sources. *Journal of*
1169         *Neuroscience*, 30(38), 12572-12580.

1170   Chen, Y. C., & Spence, C. (2010). When hearing the bark helps to identify the dog:
1171         Semantically-congruent sounds modulate the identification of masked
1172         pictures. *Cognition*, *114*(3), 389-404.

1173   Chennu, S., Noreika, V., Gueorguiev, D., Blenkmann, A., Kochen, S., Ibánez, A., ... &
1174         Bekinschtein, T. A. (2013). Expectation and attention in hierarchical auditory
1175         prediction. *Journal of Neuroscience*, *33*(27), 11194-11205.

1176   Chun, M. M., & Jiang, Y. (1998). Contextual cueingcueing: Implicit learning and memory of
1177         visual context guides spatial attention. *Cognitive psychology*, *36*(1), 28-71.

1178   Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of
1179         cognitive science. *Behavioral and brain sciences*, *36*(3), 181-204.

1180   Correa, Á., Lupiáñez, J., & Tudela, P. (2005). Attentional preparation based on temporal
1181         expectancy modulates processing at the perceptual level. *Psychonomic bulletin &*
1182         *review*, 12(2), 328-334.

1183   Coull, J. T., Frith, C. D., Büchel, C., & Nobre, A. C. (2000). Orienting attention in time:
1184         behavioural and neuroanatomical distinction between exogenous and endogenous
1185         shifts. *Neuropsychologia*, *38*(6), 808-819.

1186   Dassanayake, T. L., Michie, P. T., & Fulham, R. (2016). Effect of temporal predictability on
1187         exogenous attentional modulation of feedforward processing in the striate
1188         cortex. *International Journal of Psychophysiology*, 105, 9-16.

1189    De Meo, R., Murray, M. M., Clarke, S., & Matusz, P. J. (2015). Top-down control and early
1190        multisensory processes: chicken vs. egg. *Frontiers in integrative neuroscience*, *9*(17),
1191        1-6.

1192    Dell'Acqua, R., Sessa, P., Peressotti, F., Mulatti, C., Navarrete, E., & Grainger, J. (2010). ERP
1193        evidence for ultra-fast semantic processing in the picture–word interference
1194        paradigm. *Frontiers in psychology*, 1, 177.

1195    Desimone, R., & Duncan, J. (1995). Neural mechanisms of selective visual attention. *Annual*
1196        *Review of Neuroscience*, *18*(1), 193-222.

1197    Doehrmann, O., & Naumer, M. J. (2008). Semantics and the multisensory brain: How
1198        meaning modulates processes of audio-visual integration. *Brain Research, 1242,*
1199        136–50. https://doi.org/10.1016/J.BRAINRES.2008.03.071

1200    Duncan, J., & Humphreys, G. W. (1989). Visual search and stimulus similarity. *Psychological*
1201        *Review, 96(3),* 433–458.

1202    Eickhoff, S. B., Milham, M., & Vanderwal, T. (2020). Towards clinical applications of movie
1203        fMRI. Neuroimage, 116860.

1204    Eimer, M. (1996). The N2pc component as an indicator of attentional selectivity.
1205        *Electroencephalography and Clinical Neurophysiology, 99(3),* 225–234.

1206    Eimer, M. (2014). The neural basis of attentional control in visual search. *Trends in Cognitive*
1207        *Sciences*, *18*(10), 526-535.

1208    Eimer, M., & Kiss, M. (2008). Involuntary attentional capture is determined by task set:
1209        Evidence from event-related brain potentials. *Journal of cognitive*
1210        *neuroscience*, *20*(8), 1423-1433.

1211    Eimer, M., Kiss, M., Press, C., & Sauter, D. (2009). The roles of feature-specific task set and
1212        bottom-up salience in attentional capture: An ERP study. *Journal of Experimental*
1213        *Psychology: Human Perception and Performance, 35(5),* 1316–1328.

1214    Ernst, M. O. (2007). Learning to integrate arbitrary signals from vision and touch. *Journal of*
1215        *Vision*, *7*(5), 7-7.

1216    Fleming, J. T., Noyce, A. L., & Shinn-Cunningham, B. G. (2020). Audio-visual spatial alignment
1217        improves integration in the presence of a competing audio-visual
1218        stimulus. *Neuropsychologia*, *146*, 107530.

1219    Folk, C. L., Leber, A. B., & Egeth, H. E. (2002). Made you blink! Contingent attentional
1220        capture produces a spatial blink. *Perception & psychophysics*, *64*(5), 741-753.

1221    Folk, C. L., Remington, R. W., & Johnston, J. C. (1992). Involuntary covert orienting is
1222        contingent on attentional control settings. *Journal of Experimental Psychology:*
1223        *Human Perception and Performance, 18(4),* 1030–1044.

1224    Gaspelin, N., & Luck, S. J. (2019). Inhibition as a potential resolution to the attentional
1225        capture debate. *Current opinion in psychology*, *29*, 12-18.

1226    Gazzaley, A., & Nobre, A. C. (2012). Top-down modulation: bridging selective attention and
1227        working memory. *Trends in cognitive sciences*, 16(2), 129-135.

1228    Ghazanfar, A. A., Maier, J. X., Hoffman, K. L., & Logothetis, N. K. (2005). Multisensory
1229        integration of dynamic faces and voices in rhesus monkey auditory cortex. *Journal of*
1230        *Neuroscience*, 25(20), 5004-5012.

1231    Girelli, M., & Luck, S. J. (1997). Are the same attentional mechanisms used to detect visual
1232        search targets defined by color, orientation, and motion? *Journal of Cognitive*
1233        *Neuroscience*, *9*(2), 238-253.

1234    Golumbic, E. M. Z., Poeppel, D., & Schroeder, C. E. (2012). Temporal context in speech
1235        processing and attentional stream selection: a behavioral and neural
1236        perspective. *Brain and language*, *122*(3), 151-161.

1237    Green, J. J., & McDonald, J. J. (2010). The role of temporal predictability in the anticipatory
1238        biasing of sensory cortex during visuospatial shifts of
1239        attention. *Psychophysiology*, 47(6), 1057-1065.

1240    Grill-Spector, K., Henson, R., & Martin, A. (2006). Repetition and the brain: neural models of
1241        stimulus-specific effects. *Trends in Cognitive Sciences*, *10*(1), 14-23.

1242    Guthrie, D., & Buchwald, J. S. (1991). Significance testing of difference potentials.
1243        *Psychophysiology, 28(2),* 240-244.

1244    Hickey, C., Di Lollo, V., & McDonald, J. J. (2008). Target and distractor processing in visual
1245        search: Decomposition of the N2pc. *Visual Cognition*, *16*(1), 110-113.

1246    Hickey, C., Di Lollo, V., & McDonald, J. J. (2009). Electrophysiological indices of target and
1247        distractor processing in visual search. Journal of cognitive neuroscience, 21(4), 760-
1248        775.

1249    Holm, S. (1979). A Simple Sequentially Rejective Multiple Test Procedure. *Scandinavian*
1250        *Journal of Statistics*, *6*(2), 65-70.
1251    Hopf, J.-M., Luck, S. J., Girelli, M., Mangun, G. R., Scheich, H., & Heinze, H.-J. (2000). Neural
1252        sources of focused attention in visual search. *Cerebral Cortex, 10*, 1233–1241.
1253    Huth, A. G., Lee, T., Nishimoto, S., Bilenko, N. Y., Vu, A. T., & Gallant, J. L. (2016). Decoding
1254        the semantic content of natural movies from human brain activity. *Frontiers in*
1255        *systems neuroscience*, *10*, 81.
1256        induced gamma band responses reflect cross-modal interactions in familiar object
1257        recognition. *Journal of Neuroscience*, *27*(5), 1090-1096.
1258    Iordanescu, L., Guzman-Martinez, E., Grabowecky, M., & Suzuki, S. (2008). Characteristic
1259        sounds facilitate visual search. *Psychonomic Bulletin & Review*, *15*(3), 548-554.

1260    Jiang, J., Brashier, N. M., & Egner, T. (2015). Memory meets control in hippocampal and
1261        striatal binding of stimuli, responses, and attentional control states. *Journal of*
1262        *Neuroscience, 35*, 14885-14895.

1263    Kingstone, A., Smilek, D., Ristic, J., Kelland Friesen, C., & Eastwood, J. D. (2003). Attention,
1264        researchers! It is time to take a look at the real world. *Current Directions in*
1265        *Psychological Science*, *12*(5), 176-180.

1266    Kiss, M., Jolicœur, P., Dell'Acqua, R., & Eimer, M. (2008a). Attentional capture by visual
1267        singletons is mediated by top-down task set: New evidence from the N2pc
1268        component. *Psychophysiology*, *45*(6), 1013-1024.

1269    Kiss, M., Van Velzen, J., & Eimer, M. (2008b). The N2pc component and its links to attention
1270        shifts and spatially selective visual processing. Psychophysiology, 45, 240–249.

1271    Klaffehn, A. L., Baess, P., Kunde, W., & Pfister, R. (2019). Sensory attenuation prevails when
1272        controlling for temporal predictability of self-and externally generated
1273        tones. *Neuropsychologia*, *132*, 107145.

1274    Koenig, T., Stein, M., Grieder, M., & Kottlow, M. (2014). A tutorial on data-driven methods
1275        for statistically assessing ERP topographies. *Brain topography*, *27*(1), 72-83.

1276    Kuo, B. C., Nobre, A. C., Scerif, G., & Astle, D. E. (2016). Top–Down activation of spatiotopic
1277        sensory codes in perceptual and working memory search. *Journal of cognitive*
1278        *neuroscience*, *28*(7), 996-1009.

1279    Laurienti, P. J., Burdette, J. H., Maldjian, J. A., & Wallace, M. T. (2006). Enhanced
1280        multisensory integration in older adults. *Neurobiology of aging*, *27*(8), 1155-1163.

1281    Lehmann, D. (1987). "Principles of spatial analysis," in *Methods of Analysis of Brain Electrical*
1282        *and Magnetic Signals*, A. S. Gevins and A. Remont, Eds., pp. 309–354. Elsevier:
1283        Amsterdam, The Netherlands.

1284    Lehmann D, Skrandies W (1980): Reference-free identification of components of
1285        checkerboard evoked multichannel potential fields. *Electroencephalography in*
1286        *Clinical Neurology, 48*, 609–621.

1287    Lehmann, D., Ozaki, H., & Pal, I. (1987). EEG alpha map series: brain micro-states by space-
1288        oriented adaptive segmentation. *Electroencephalography and clinical*
1289        *neurophysiology*, *67*(3), 271-288.

1290    Lien, M. C., Ruthruff, E., Goodin, Z., & Remington, R. W. (2008). Contingent attentional
1291        capture by top-down control settings: converging evidence from event-related
1292        potentials. *Journal of Experimental Psychology: Human Perception and*
1293        *Performance*, *34*(3), 509.

1294    Luck, S. J., Gaspelin, N., Folk, C. L., Remington, R. W., & Theeuwes, J. (2020). Progress
1295        toward resolving the attentional capture debate. *Visual Cognition*, 1-21. DOI:
1296        10.1080/13506285.2020.1848949

1297    Luck, S. J., & Hillyard, S. A. (1994a). Electrophysiological correlates of feature analysis during
1298        visual search. *Psychophysiology, 31,* 291–308.

1299    Luck, S. J., & Hillyard, S. A. (1994b). Spatial filtering during visual search: Evidence from
1300        human electrophysiology. *Journal of Experimental Psychology: Human Perception*
1301        *and Performance, 20(5),* 1000–1014.

1302    Lunn, J., Sjoblom, A., Ward, J., Soto-Faraco, S., & Forster, S. (2019). Multisensory
1303        enhancement of attention depends on whether you are already paying
1304        attention. *Cognition*, *187*, 38-49.

1305    Luo, H., & Poeppel, D. (2007). Phase patterns of neuronal responses reliably discriminate
1306        speech in human auditory cortex. *Neuron*, *54*(6), 1001-1010.

1307    Matusz, P. J., & Eimer, M. (2013). Top-down control of audiovisual search by bimodal search
1308        templates. *Psychophysiology*, *50*(10), 996-1009.

1309    Matusz, P. J., & Eimer, M. (2011). Multisensory enhancement of attentional capture in visual
1310        search. *Psychonomic bulletin & review*, *18*(5), 904.

1311    Matusz, P. J., Wallace, M. T., & Murray, M. M. (2020). Multisensory contributions to object
1312        recognition and memory across the life span. In *Multisensory Perception* (pp. 135-
1313        154). Academic Press.

1314    Matusz, P. J., Dikker, S., Huth, A. G., & Perrodin, C. (2019a). Are We Ready for Real-world
1315        Neuroscience?. *Journal of cognitive neuroscience*, *31*(3), 327.

1316   Matusz, P. J., Turoman, N., Tivadar, R. I., Retsa, C., & Murray, M. M. (2019b). Brain and
1317        cognitive mechanisms of top–down attentional control in a multisensory world:
1318        Benefits of electrical neuroimaging. *Journal of cognitive neuroscience*, *31*(3), 412-
1319        430.

1320   Matusz, P. J., Merkley, R., Faure, M., & Scerif, G. (2019c). Expert attention: Attentional
1321        allocation depends on the differential development of multisensory number
1322        representations. *Cognition*, *186*, 171-177.

1323   Matusz, P. J., Key, A. P., Gogliotti, S., Pearson, J., Auld, M. L., Murray, M. M., & Maitre, N. L.
1324        (2018). Somatosensory plasticity in pediatric cerebral palsy following constraint-
1325        induced movement therapy. *Neural Plasticity*, *2018*.

1326   Matusz, P. J., Wallace, M. T., & Murray, M. M. (2017). A multisensory perspective on object
1327        memory. *Neuropsychologia*, *105*, 243-252.

1328   Matusz, P. J., Retsa, C., & Murray, M. M. (2016). The context-contingent nature of cross-
1329        modal activations of the visual cortex. *Neuroimage*, *125*, 996-1004.

1330   Matusz, P. J., Thelen, A., Amrein, S., Geiser, E., Anken, J., & Murray, M. M. (2015a). The role
1331        of auditory cortices in the retrieval of single-trial auditory–visual object
1332        memories. *European Journal of Neuroscience*, *41*(5), 699-708.

1333   Matusz, P. J., Broadbent, H., Ferrari, J., Forrest, B., Merkley, R., & Scerif, G. (2015b). Multi-
1334        modal distraction: Insights from children's limited attention. *Cognition*, *136*, 156-
1335        165.

1336   McNab, F., & Klingberg, T. (2008). Prefrontal cortex and basal ganglia control access to
1337        working memory. *Nature Neuroscience*, *11*, 103-107.

1338   Michel, C. M., & Murray, M. M. (2012). Towards the utilization of EEG as a brain imaging
1339        tool. *Neuroimage*, *61*(2), 371-385.

1340   Miniussi C, Wilding EL, Coull JT, Nobre AC. (1999). Orienting atten- tion in the time domain:
1341        modulation of potentials. *Brain, 122,* 1507-18.

1342   Murray, M. M., Brunet, D., & Michel, C. M. (2008). Topographic ERP analyses: a step-by-step
1343        tutorial review. *Brain topography*, *20*(4), 249-264.

1344   Murray, M. M., Thelen, A., Thut, G., Romei, V., Martuzzi, R., & Matusz, P. J. (2016a). The
1345        multisensory function of the human primary visual cortex. *Neuropsychologia*, *83*,
1346        161-169.

1347   Murray, M. M., Lewkowicz, D. J., Amedi, A., & Wallace, M. T. (2016b). Multisensory
1348        processes: a balancing act across the lifespan. *Trends in Neurosciences*, *39*(8), 567-
1349        579.

1350   Murray, M. M., Michel, C. M., De Peralta, R. G., Ortigue, S., Brunet, D., Andino, S. G., &
1351        Schnider, A. (2004). Rapid discrimination of visual and multisensory memories
1352        revealed by electrical neuroimaging. *Neuroimage*, *21*(1), 125-135.

1353   Naccache, L., Blandin, E., & Dehaene, S. (2002). Unconscious masked priming depends on
1354        temporal attention. *Psychological Science*, 13(5), 416-424.

1355   Nastase, S. A., Goldstein, A., & Hasson, U. (2020). Keep it real: rethinking the primacy of
1356        experimental control in cognitive neuroscience. *Neuroimage*. In print.

1357    Naumann, S., Byrne, M. L., de la Fuente, L. A., Harrewijn, A., Nugiel, T., Rosen, M. L., ... &
1358        Matusz, P. J. (2020). Assessing the degree of ecological validity of your study:
1359        Introducing the Ecological Validity Assessment (EVA) Tool. *PsyArXiv*. DOI:
1360        10.31234/osf.io/qb9tz.

1361    Neel, M. L., Yoder, P., Matusz, P. J., Murray, M. M., Miller, A., Burkhardt, S., ... & Maitre, N.
1362        L. (2019). Randomized controlled trial protocol to improve multisensory neural
1363        processing, language and motor outcomes in preterm infants. *BMC Pediatrics*, *19*(1),
1364        1-10.

1365    Noonan, M. P., Crittenden, B. M., Jensen, O., & Stokes, M. G. (2018). Selective inhibition of
1366        distracting input. *Behavioural brain research*, *355*, 36-47.

1367    Peelen, M. V., & Kastner, S. (2014). Attention in the real world: toward understanding its
1368        neural basis. *Trends in cognitive sciences*, *18*(5), 242-250.

1369    Perrin, F., Pernier, J., Bertnard, O., Giard, M. H., & Echallier, J. F. (1987). Mapping of scalp
1370        potentials by surface spline interpolation. *Electroencephalography and clinical
1371        neurophysiology*, *66*(1), 75-81.

1372    Press, C., Kok, P., & Yon, D. (2020). The perceptual prediction paradox. *Trends in Cognitive
1373        Sciences*, *24*(1), 13-24.

1374    Raij, T., Ahveninen, J., Lin, F. H., Witzel, T., Jääskeläinen, I. P., Letham, B., ... & Hämäläinen,
1375        M. (2010). Onset timing of cross-sensory activations and multisensory interactions in
1376        auditory and visual sensory cortices. *European Journal of Neuroscience*, 31(10),
1377        1772-1782.

1378    Raij, T., Uutela, K., & Hari, R. (2000). Audiovisual integration of letters in the human
1379        brain. *Neuron*, *28*(2), 617-625.

1380    Retsa, C., Matusz, P. J., Schnupp, J. W., & Murray, M. M. (2018). What's what in auditory
1381        cortices?. *NeuroImage*, *176*, 29-40.

1382    Retsa, C., Matusz, P. J., Schnupp, J. W., & Murray, M. M. (2020). Selective attention to sound
1383        features mediates cross-modal activation of visual cortices. *Neuropsychologia*, *144*,
1384        107498.

1385    Richter, D., Ekman, M., & de Lange, F. P. (2018). Suppressed sensory response to predictable
1386        object stimuli throughout the ventral visual stream. *Journal of Neuroscience*, *38*(34),
1387        7452-7461.

1388    Rohenkohl, G., Gould, I. C., Pessoa, J., & Nobre, A. C. (2014). Combining spatial and temporal
1389        expectations to improve visual perception. *Journal of vision*, 14(4), 8-8.

1390    Sadeh, T., Shohamy, D., Levy, D. R., Reggev, N., & Maril, A. (2011). Cooperation between the
1391        hippocampus and the striatum during episodic encoding. Journal of Cognitive
1392        Neuroscience, 23(7), 1597-1608.

1393    Sarmiento, B. R., Matusz, P. J., Sanabria, D., & Murray, M. M. (2016). Contextual factors
1394        multiplex to control multisensory processes. Human brain mapping, 37(1), 273-288.

1395    Sawaki, R., & Luck, S. J. (2010). Capture versus suppression of attention by salient
1396        singletons: Electrophysiological evidence for an automatic attend-to-me
1397        signal. *Attention, Perception, & Psychophysics*, *72*(6), 1455-1470.

1398 Schröger, E., Marzecová, A., & SanMiguel, I. (2015). Attention and prediction in human
1399      audition: A lesson from cognitive psychophysiology. *European Journal of*
1400      *Neuroscience, 41*(5), 641-664.

1401 Shamay-Tsoory, S. G., & Mendelsohn, A. (2019). Real-life neuroscience: an ecological
1402      approach to brain and behavior research. *Perspectives on Psychological*
1403      *Science*, *14*(5), 841-859.

1404 Soto-Faraco, S., Kvasova, D., Biau, E., Ikumi, N., Ruzzoli, M., Morís-Fernández, L., & Torralba,
1405      M. (2019). *Multisensory interactions in the real world*. Cambridge University Press.

1406 Southwell, R., Baumann, A., Gal, C., Barascud, N., Friston, K., & Chait, M. (2017). Is
1407      predictability salient? A study of attentional capture by auditory
1408      patterns. *Philosophical Transactions of the Royal Society B: Biological*
1409      *Sciences*, 372(1714), 20160105.

1410 Spierer, L., Manuel, A. L., Bueti, D., & Murray, M. M. (2013). Contributions of pitch and
1411      bandwidth to sound-induced enhancement of visual cortex excitability in
1412      humans. *Cortex*, *49*(10), 2728-2734.

1413 Sui, J., He, X., & Humphreys, G. W. (2012). Perceptual effects of social salience: evidence
1414      from self-prioritization effects on perceptual matching. *Journal of Experimental*
1415      *Psychology: Human perception and performance*, *38*(5), 1105.

1416 Summerfield, C., & Egner, T. (2009). Expectation (and attention) in visual cognition. *Trends*
1417      *in cognitive sciences*, *13*(9), 403-409.

1418 Summerfield, J. J., Lepsien, J., Gitelman, D. R., Mesulam, M. M., & Nobre, A. C. (2006).
1419      Orienting attention based on long-term memory experience. *Neuron*, *49*(6), 905-916.

1420 Sun, Y., Fuentes, L. J., Humphreys, G. W., & Sui, J. (2016). Try to see it my way: Embodied
1421      perspective enhances self and friend-biases in perceptual matching. *Cognition*, *153*,
1422      108-117.

1423 Talsma, D., & Woldorff, M. G. (2005). Selective attention and multisensory integration:
1424      multiple phases of effects on the evoked brain activity. *Journal of Cognitive*
1425      *Neuroscience,* 17, 1098–1114.

1426 Ten Oever, S., & Sack, A. T. (2015). Oscillatory phase shapes syllable perception. *Proceedings*
1427      *of the National Academy of Sciences*, *112*(52), 15833-15837.

1428 Ten Oever, S., Romei, V., van Atteveldt, N., Soto-Faraco, S., Murray, M. M., & Matusz, P. J.
1429      (2016). The COGs (context, object, and goals) in multisensory
1430      processing. *Experimental brain research*, 234(5), 1307-1323.

1431 Theeuwes, J. (1991). Cross-dimensional perceptual selectivity. *Perception & Psychophysics,*
1432      *50(2),* 184–193.

1433 Thelen, A., Talsma, D., & Murray, M.M. (2015). Single-trial multisensory memories affect
1434      later auditory and visual object discrimination. *Cognition 138,* 148–160.

1435 Thorpe, S., Fize, D., & Marlot, C. (1996). Speed of processing in the human visual
1436      system. *Nature*, *381*(6582), 520-522.

1437 Tivadar, R. I., & Murray, M. M. (2019). A primer on electroencephalography and event-
1438      related potentials for organizational neuroscience. *Organizational Research*
1439      *Methods*, *22*(1), 69-94.

1440   Tivadar, R. I., Knight, R. T., & Tzovara, A. (2021). Automatic Sensory Predictions: A Review of
1441        Predictive Mechanisms in the Brain and Their Link to Conscious Processing. *Frontiers*
1442        *in Human Neuroscience*, 438.

1443   Tovar, D. A., Murray, M. M., & Wallace, M. T. (2020). Selective enhancement of object
1444        representations through multisensory integration. *Journal of Neuroscience*. In press.
1445        DOI: https://doi.org/10.1523/JNEUROSCI.2139-19.2020

1446   Treisman, A. M., & Gelade, G. (1980). A feature-integration theory of attention. *Cognitive*
1447        *Psychology, 12(1),* 97–136.

1448   Turoman, N., Tivadar, R. I., Retsa, C., Maillard, A. M., Scerif, G., and Matusz, P. J. (2021a).
1449        The development of attentional control mechanisms in multisensory environments.
1450        *Developmental Cognitive Neuroscience, 48,* 100930.

1451   Turoman, N., Tivadar, R. I., Retsa, C., Maillard, A. M., Scerif, G., and Matusz, P. (2021b).
1452        Uncovering the mechanisms of real-world attentional control over the course of
1453        primary education. *Mind, Brain, & Education.* In press.

1454   Tzovara, A., Murray, M. M., Michel, C. M., & De Lucia, M. (2012). A tutorial review of
1455        electrical   neuroimaging   from   group-average   to   single-trial   event-related
1456        potentials. *Developmental neuropsychology*, *37*(6), 518-544.

1457   Van Atteveldt, N., Murray, M. M., Thut, G., & Schroeder, C. E. (2014). Multisensory
1458        integration: flexible use of general operations. *Neuron*, *81*(6), 1240-1253.

1459   van Atteveldt, N., van Kesteren, M. T. R., Braams, B.; Krabbendam, L. (2018). Neuroimaging
1460        of   learning   and   development:   improving   ecological   validity. *Frontline Learning*
1461        *Research,* 6 (3), 186–203. DOI: 10.14786/flr.v6i3.366.

1462   Van der Burg, E., Talsma, D., Olivers, C. N. L., Hickey, C., & Theeuwes, J. (2011). Early
1463        multisensory   interactions   affect   the   competition   among   multiple   visual   objects.
1464        *NeuroImage, 55,* 1208–1218.

1465   van Moorselaar, D., & Slagter, H. A. (2019). Learning what is irrelevant or relevant:
1466        Expectations facilitate distractor inhibition and target facilitation through distinct
1467        neural mechanisms. *Journal of Neuroscience*, *39*(35), 6953-6967.

1468   van Moorselaar, D., & Slagter, H. A. (2020). Inhibition in selective attention. *Annals of the*
1469        *New York Academy of Sciences*, *1464*(1), 204.

1470   van Moorselaar, D., Daneshtalab, N., & Slagter, H. (2020). Neural mechanisms underlying
1471        distractor inhibition on the basis of feature and/or spatial expectations. bioRxiv.

1472   Vanderwal, T., Eilbott, J.; Castellanos, F. X (2019). Movies in the magnet: Naturalistic
1473        paradigms   in   developmental   functional   neuroimaging. *Developmental Cognitive*
1474        *Neuroscience, 36*, 100600.

1475   Vaughan Jr, H. G. (1982). The neural origins of human event-related potentials. *Annals of the*
1476        *New York Academy of Sciences*, *388*(1), 125-138.

1477   Widmann, A., Schröger, E., & Maess, B. (2015). Digital filter design for electrophysiological
1478        data–a practical approach. *Journal of Neuroscience Methods, 250*, 34-46.

1479   Wu, R., Nako, R., Band, J., Pizzuto, J., Ghoreishi, Y., Scerif, G., & Aslin, R. (2015). Rapid
1480        attentional selection of non-native stimuli despite perceptual narrowing. *Journal of*
1481        *Cognitive Neuroscience, 27*(11), 2299-2307.

1482  Yuval-Greenberg, S., & Deouell, L. Y. (2007). What you see is not (always) what you hear:
1483      induced gamma band responses reflect cross-modal interactions in familiar object
1484      recognition. *Journal of Neuroscience*, *27*(5), 1090-1096.

1485
1486
1487
1488
1489
1490
1491
1492
1493
1494  **Acknowledgments**
1495
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528

1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544 **Figure Legends**
1545
1546 **Figure 1. A)** An example trial of the general experimental "Task" is shown, with four
1547 successive arrays. The white circle around the target location (here the target is a blue
1548 diamond) and the corresponding distractor location serves to highlight, in this case, a target-
1549 matching distractor colour condition, with a concomitant sound, i.e., TCCAV. **B)** The order of
1550 Tasks, with the corresponding conditions of Multisensory Relationship in red, and Distractor
1551 Onset in green, shown separately for each Task, in the successive order in which they
1552 appeared in the study. Under each condition, its operationalisation is given in brackets in
1553 the corresponding colour. Predictable and unpredictable blocks before and after the
1554 training (1 & 2 and 3 & 4, respectively) were counterbalanced across participants. **C)** Events
1555 that were part of the Training. Association phase: an example pairing option (red – high
1556 pitch, blue – low pitch) with trial progression is shown. Testing phase: the pairing learnt in
1557 the Association phase would be tested using a colour word or a string of x's in the respective
1558 colour. Participants had to indicate whether the pairing was correct via a button press, after
1559 which feedback was given.
1560
1561 **Figure 2.** The violin plots show the attentional capture effects (spatial cueing in
1562 milliseconds) for TCC and NCC distractors, and the distributions of single-participant scores
1563 according to whether Multisensory Relationship within these distractors was Arbitrary (light
1564 green) or Congruent (dark green). The dark grey boxes within each violin plot show the
1565 interquartile range from the $1^{st}$ to the $3^{rd}$ quartile, and white dots in the middle of these
1566 boxes represent the median. Larger values indicate *positive* behavioural capture effects (RTs
1567 faster on trials where distractor and target appeared in same vs. different location), while
1568 below-zero values – *inverted* capture effects (RTs slower on trials where distractor and
1569 target appeared in same vs. different location). Larger behavioural capture elicited by
1570 target-colour distractors (TCC) was found for arbitrary than semantically congruent
1571 distractors. Expectedly, regardless of Multisensory Relationship, attentional capture was
1572 larger for target-colour (TCC) distractors than for non-target colour distractors (NCC).
1573
1574 **Figure 3.** Overall contra- and ipsilateral ERP waveforms representing a mean amplitude over
1575 electrode clusters (plotted on the head model at the bottom of the figure in blue and black),

36

1576    separately for each of the four experimental conditions (Distractor Colour x Distractor
1577    Modality), averaged across all four Tasks. The N2pc time-window of 180–300ms following
1578    distractor onset is highlighted in grey, and significant contra-ipsi differences are marked
1579    with an asterisk ($p < 0.05$). As expected, only the TCC distractors elicited statistically
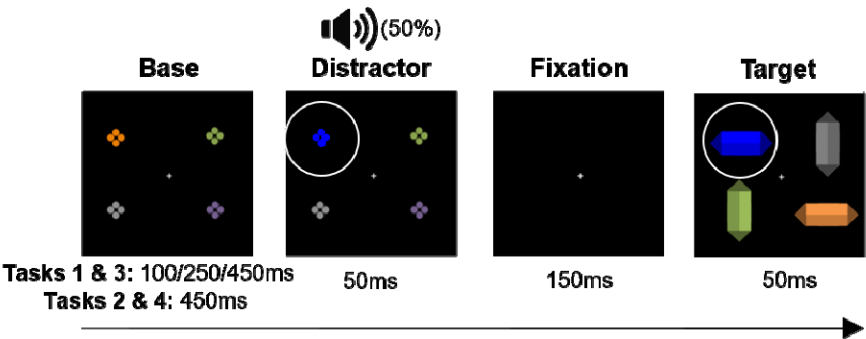1580    significant contra-ipsi differences.
1581
1582    **Figure 4.** Nonlateralised GFP and topography results for the visual only difference ERPs (DV
1583    condition of Target Difference), as a proxy for TAC. **A)** Mean GFP over the post-distractor
1584    and pre-target time-period across the 4 experimental tasks (as a function of the levels of
1585    Multisensory Relationship and Distractor Onset that they represent), as denoted by the
1586    colours on the legend. The time-windows of interest (102–124ms and 234–249ms) are
1587    highlighted by grey areas. **B)** Template maps over the post-distractor time-period as
1588    revealed by the topographic clustering (Maps A1 to A5) are shown in top panels. In lower
1589    panels are the results of the fitting procedure over the 29–126ms time-window. The results
1590    displayed here are the follow-up tests of the 3-way Map x Multisensory Relationship x
1591    Distractor Onset interaction as a function of Multisensory Relationship (leftward panel) and
1592    of Distractor Onset (rightward panel). Bars are coloured according to the template maps
1593    that they represent. Conditions are represented by full colour or patterns per the legend.
1594    Error bars represent standard errors of the mean.
1595
1596    **Figure 5.** Nonlateralised GFP and topography results for the difference ERPs between the
1597    DAV and DV conditions of Target Difference, as a proxy for MSE. **A)** Mean GFP over the post-
1598    distractor and pre-target time-period across the 4 experimental tasks (as a function of the
1599    levels of Multisensory Relationship and Distractor Onset that they represent), as denoted by
1600    the colours on the legend. The time-windows of interest (102–124ms and 234–249ms) are
1601    highlighted by grey bars. **B)** Template maps over the post-distractor time-period as revealed
1602    by the topographic clustering (Maps A1 to A7) are shown on top. Below are the results of
1603    the fitting procedure over the three time-windows: 35–110, 110–190, and 190–300 time-
1604    window. Here we display the follow-ups of the interactions observed in each time-window:
1605    in 35–110 and 190–300 time-windows, the 2-way Map x Multisensory Relationship
1606    interaction (leftward and rightward panels, respectively), and in the 110–190 time-window,
1607    follow-ups of the 3-way Map x Multisensory Relationship x Distractor Onset interaction as a
1608    function of Multisensory Relationship and of Distractor Onset (middle panel). Bars are
1609    coloured according to the template maps that they represent. Conditions are represented
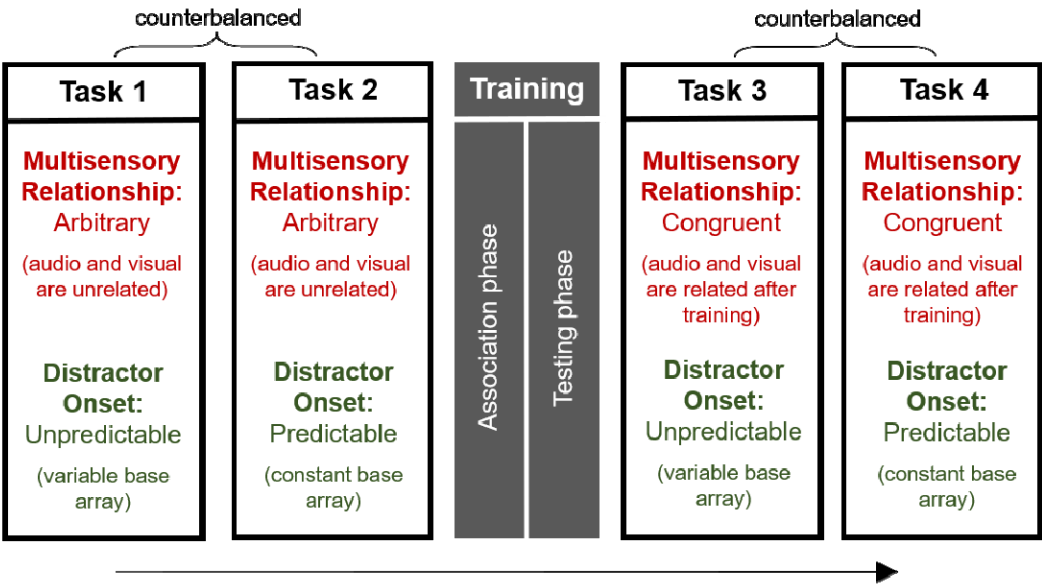1610    by full colour or patterns per the legend. Error bars represent standard errors of the mean.
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622

1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
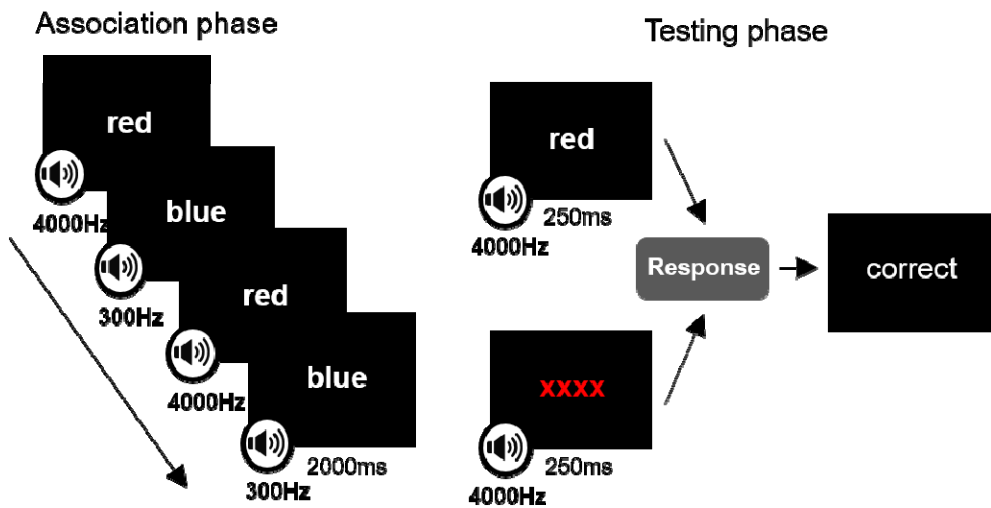1634
1635
1636
1637
1638    **Figure 1**

## A) General trial sequence across Tasks



## B) Overall structure of the study



## C) Training of semantic audio-visual associations for distractors



1639
1640
1641    **Figure 2**

# Behavioural attentional capture

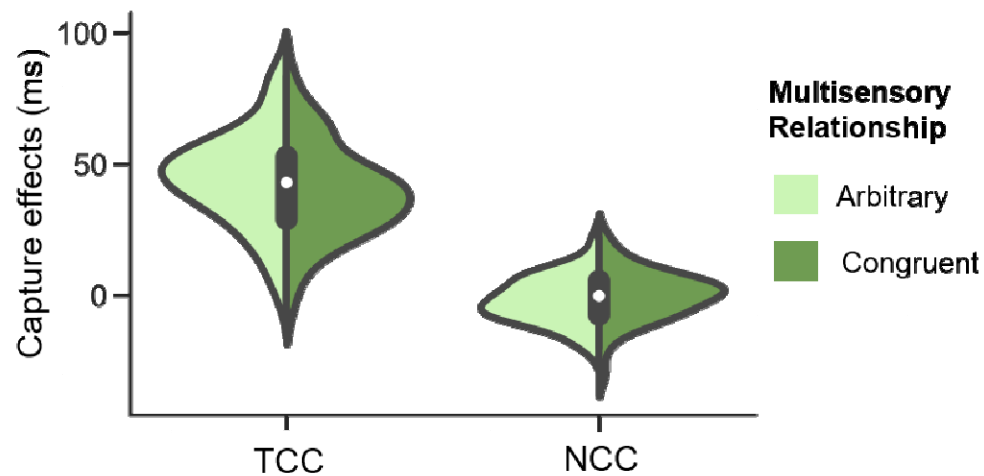## Interaction between Task-set contingent attentional capture and Multisensory Relationship



1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
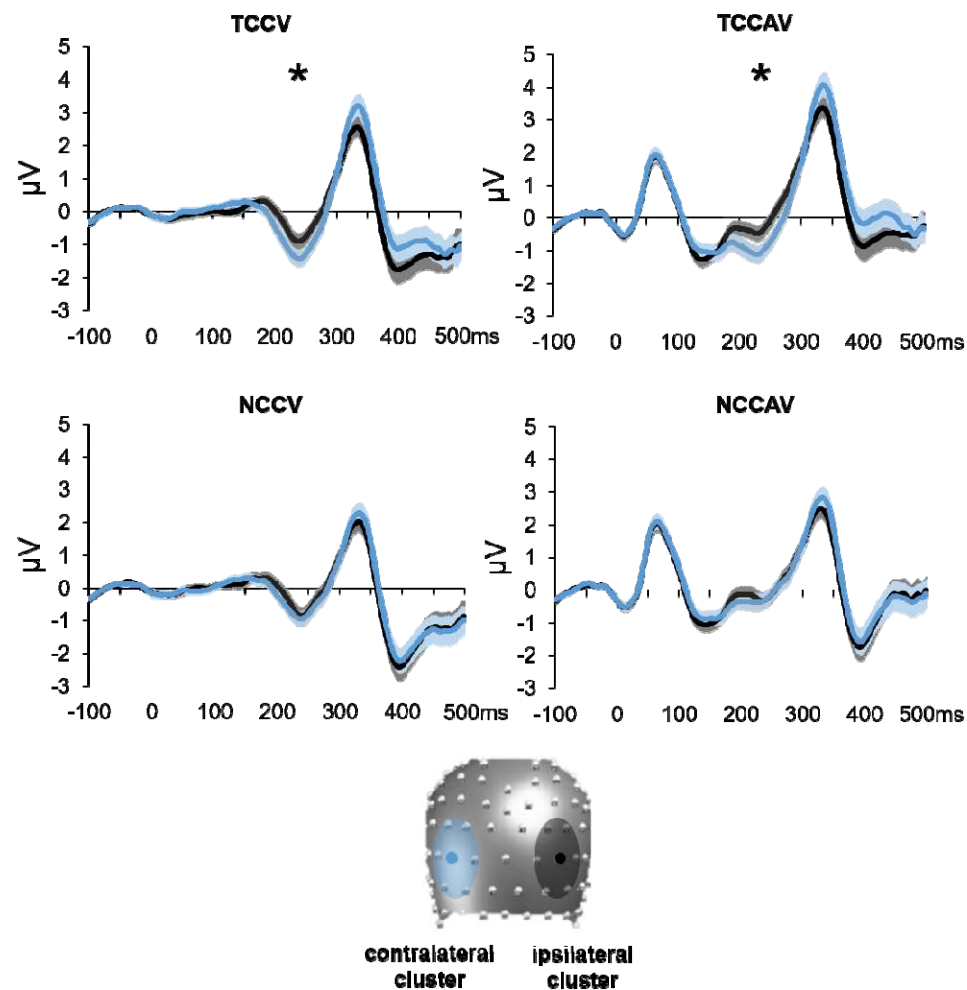1661
1662
1663
1664
1665
1666
1667
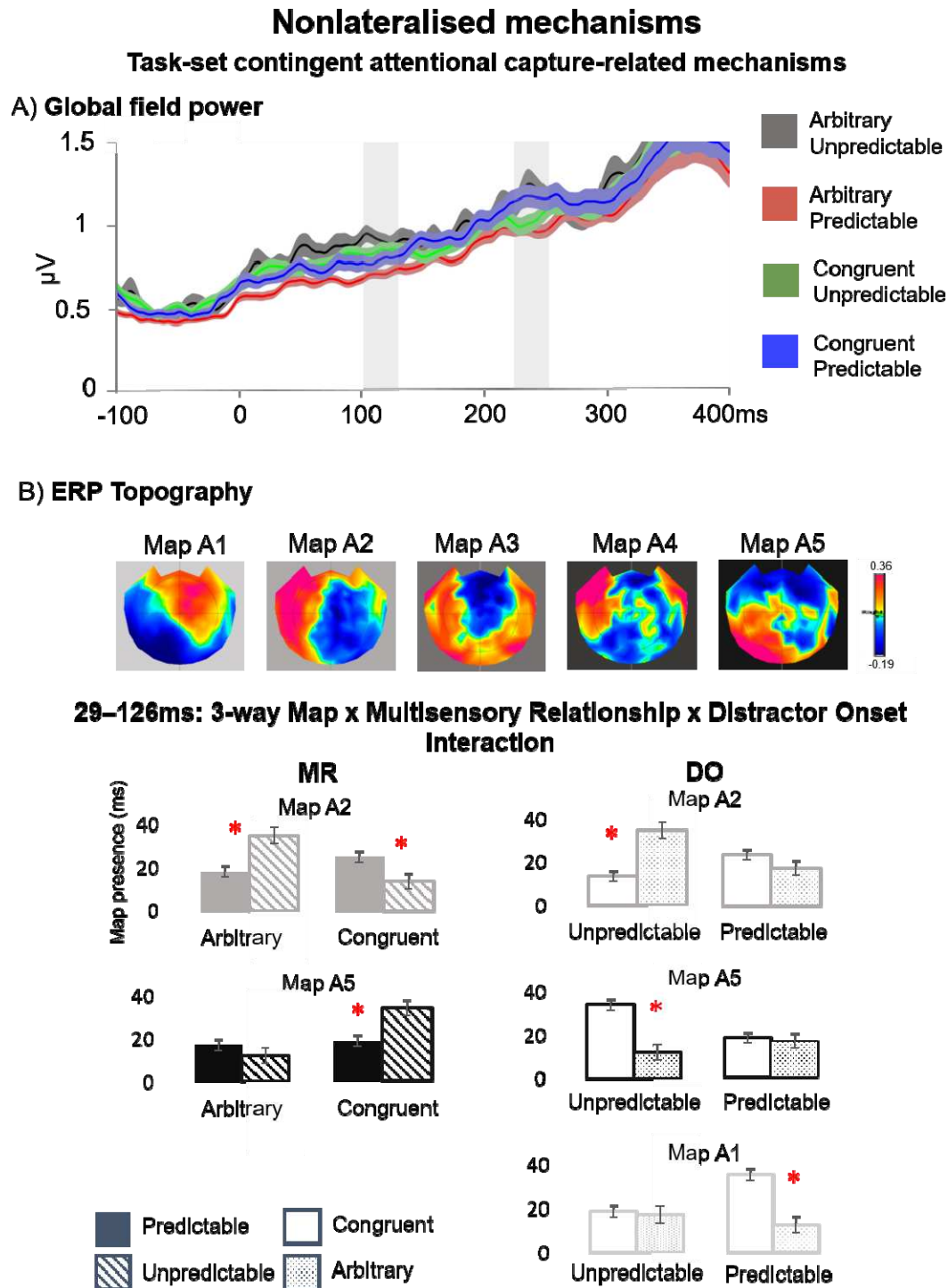1668
1669      **Figure 3**

1670

## Contralateral–Ipsilateral waveforms across experiments



1671
1672
1673
1674
1675
1676
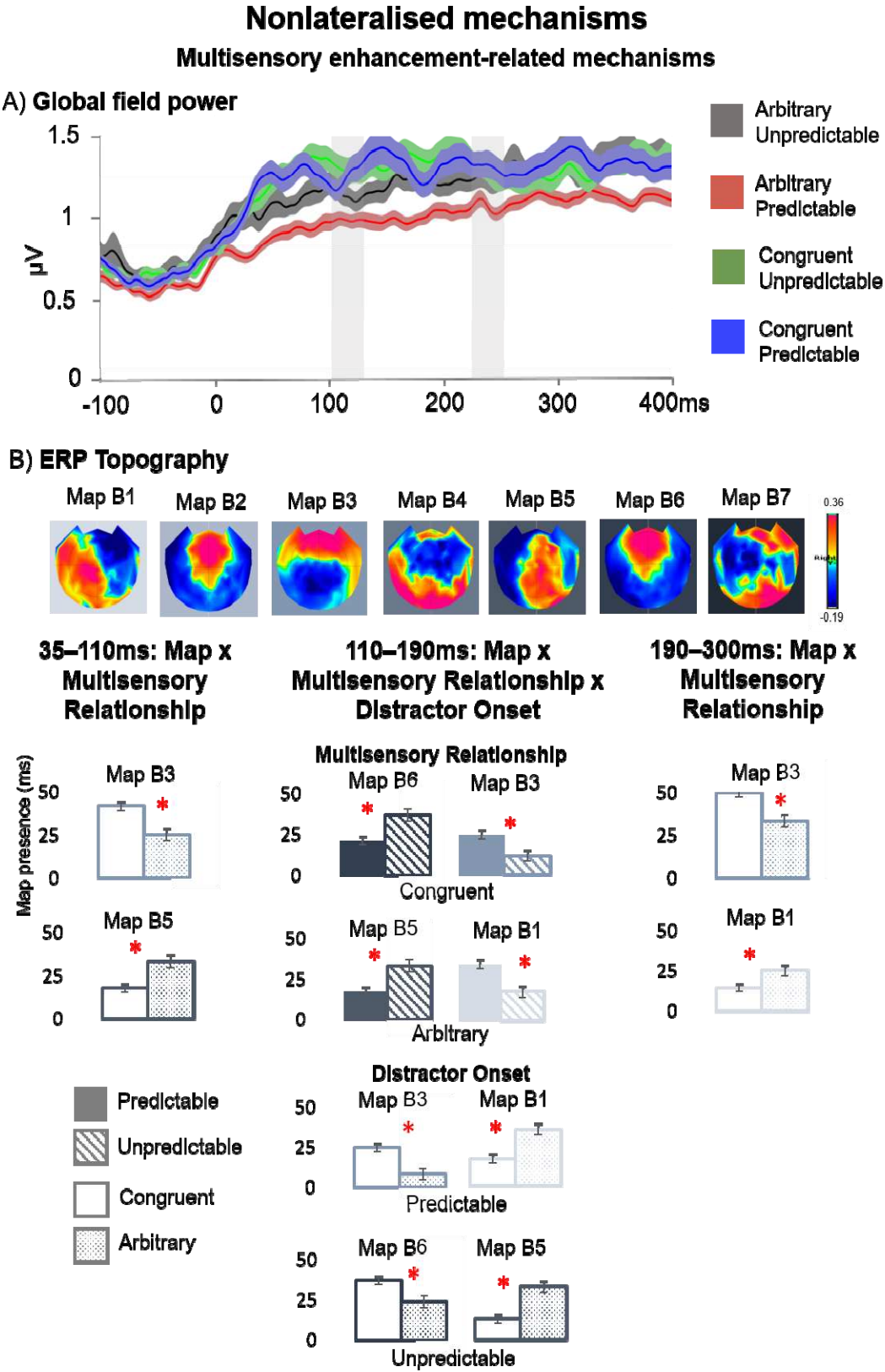1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688 **Figure 4**

1689
1690
1691



# Nonlateralised mechanisms
## Task-set contingent attentional capture-related mechanisms

1692
1693
1694
1695
1696 **Figure 5**

# Nonlateralised mechanisms

## Multisensory enhancement-related mechanisms



A) **Global field power**

B) **ERP Topography**

1697

**Appendix 1. Abbreviations**

1699

1700 N2pc – the N2pc event-related component

1701 EEG – Electroencephalography

1702 ERPs – Event-Related Potentials

1703 TAC – Task-set Contingent Attentional Capture

1704 MSE – Multisensory Enhancement of Attentional Capture

1705 SOMs – Supplementary Online Materials

1706 TCCV – target-color cue visual

1707 NCCV – nontarget-color cue visual

1708 TCCAV – target-color cue audiovisual

1709 NCCAV – nontarget-color cue audiovisual

1710 rmANOVA – repeated-measures analysis of variance

1711 GFP – Global Field Power

1712 TAAHC – Topographic Atomize and Agglomerate Hierarchical Clustering

1713 $D_{AV}$ – Target Difference, difference between TCCAV and NCCAV conditions

1714 $D_{V}$ – Target Difference, difference between TCCV and NCCV conditions

1715 DO – Distractor Onset

1716 MR – Multisensory Relationship