

Dear Author

Here are the proofs of your article.

- You can submit your corrections **online** or by **fax**.
- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.
- Please return your proof together with the permission to publish confirmation.
- For **fax** submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.
- Remember to note the journal title, article number, and your name when sending your response via e-mail, fax or regular mail.
- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.
- **Check** the questions that may have arisen during copy editing and insert your answers/corrections.
- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.
- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.
- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style. Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.
- If we do not receive your corrections **within 48 hours**, we will send you a reminder.

Please note

Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI.

Further changes are, therefore, not possible.

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL:

<http://dx.doi.org/10.3758/s13414-020-02192-y>

If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information, go to:

<http://www.springerlink.com>.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us, if you would like to have these documents returned.

The **printed version** will follow in a forthcoming issue.

Metadata of the article that will be visualized in OnlineFirst

1	Article Title	Ensemble perception: Extracting the average of perceptual versus numerical stimuli	
2	Article Sub- Title		
3	Article Copyright - Year	The Psychonomic Society, Inc. 2020 (This will be the copyright line in the final PDF)	
4	Journal Name	Attention, Perception, & Psychophysics	
5		Family Name	Rosenbaum
6		Particle	
7		Given Name	David
8	Corresponding	Suffix	
9	Author	Organization	Tel Aviv University
10		Division	
11		Address	Tel Aviv-Yafo, Israel
12		e-mail	davidros28@gmail.com
13		Family Name	Gardelle
14		Particle	de
15		Given Name	Vincent
16		Suffix	
17	Author	Organization	Paris School of Economics and CNRS, Centre d'Economie de la Sorbonne
18		Division	
19		Address	Paris, France
20		e-mail	
21		Family Name	Usher
22		Particle	
23		Given Name	Marius
24		Suffix	
25	Author	Organization	Tel Aviv University
26		Division	
27		Address	Tel Aviv-Yafo, Israel
28		e-mail	
29	Schedule	Received	
30		Revised	

31		Accepted	25 October 2020
32	Abstract	<p>Recent research has established that humans can extract an average perceptual feature over briefly presented arrays of visual elements or the average of a rapid temporal sequence of numbers. Here we compared the extraction of the average over briefly presented arrays, for a perceptual feature (orientations) and for numerical values (1–9 digits), using an identical experimental design for the two tasks. We hypothesized that the averaging of numbers, more than of orientations, would be constrained by capacity limitations. Arrays of Gabor elements or digits were simultaneously presented for 300 ms and observers were required to estimate the average on a continuous response scale. In each trial the elements were sampled from normal distributions (of various means) and we varied the set size (4–12). We found that while for orientation the averaging precision remained constant with set size, for numbers it decreased with set size. Using computational modeling we also extracted capacity parameters (the number of elements that are pooled in the average extraction). Despite marked heterogeneity between observers, the capacity for orientations (around eight items) was much larger than for numbers (around four items). The orientation task also had a larger fraction of participants relying on distributed attention to all elements. Our study thus supports the idea that numbers more than perceptual features are subject to capacity or attentional limitations when observers need to evaluate the average over an ensemble of stimuli.</p>	
33	Keywords separated by ' - '	Cognitive neuroscience - Decision making - Math modeling	
34	Foot note information	Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.	



Ensemble perception: Extracting the average of perceptual versus numerical stimuli

David Rosenbaum¹ · Vincent de Gardelle² · Marius Usher¹Accepted: 25 October 2020
© The Psychonomic Society, Inc. 2020

Abstract

Recent research has established that humans can extract an average perceptual feature over briefly presented arrays of visual elements or the average of a rapid temporal sequence of numbers. Here we compared the extraction of the average over briefly presented arrays, for a perceptual feature (orientations) and for numerical values (1–9 digits), using an identical experimental design for the two tasks. We hypothesized that the averaging of numbers, more than of orientations, would be constrained by capacity limitations. Arrays of Gabor elements or digits were simultaneously presented for 300 ms and observers were required to estimate the average on a continuous response scale. In each trial the elements were sampled from normal distributions (of various means) and we varied the set size (4–12). We found that while for orientation the averaging precision remained constant with set size, for numbers it decreased with set size. Using computational modeling we also extracted capacity parameters (the number of elements that are pooled in the average extraction). Despite marked heterogeneity between observers, the capacity for orientations (around eight items) was much larger than for numbers (around four items). The orientation task also had a larger fraction of participants relying on distributed attention to all elements. Our study thus supports the idea that numbers more than perceptual features are subject to capacity or attentional limitations when observers need to evaluate the average over an ensemble of stimuli.

Keywords Cognitive neuroscience · Decision making · Math modeling

Introduction

Research over the last two decades indicates that human observers can rapidly extract the average of a perceptual feature over sets of visual objects, even when they cannot discriminate if an individual item in the display was presented (Ariely, 2001; Chong & Treisman, 2003; Chong & Treisman, 2005; Dakin, 2001; Parkes, Lund, Angelucci, Solomon, & Morgan, 2001; Robitaille & Harris, 2011). For example, humans can evaluate the average size of a set of circles presented simultaneously, with an accuracy that does not decrease as the set contains more elements (Ariely, 2001; Chong & Treisman, 2005) or is presented for a shorter duration (Chong & Treisman, 2003). This averaging ability has been demonstrated even in situations where the discrimination of the presence of individual elements in the array appears at chance (Ariely, 2001). This capacity appears to

extend from simple visual attributes – such as size, orientation, and spatial position – to more complex properties such as emotional expression (Haberman & Whitney, 2011). Moreover, the extraction of the average appears to take place automatically or, at least, without “intention,” as it occurs in parallel (Chong & Treisman, 2005) and affects judgments of memory, in which the set-average is task-irrelevant (Khayat & Hochstein, 2018).

Another type of stimulus in which ensemble perception has been suggested is symbolic numbers (Brezis, Bronfman, Jacoby, Lavidor, & Usher, 2016; Brezis et al., 2015, 2018; Corbett, Oriet, & Rensink, 2006; Vanunu, Hotaling, & Newell, 2020; Sato & Motoyoshi, 2020; Van Opstal et al., 2011; Vandormael, Herce, Balaguer, Li, & Summerfield (2017); Spietzer et al., 2017). Such stimuli are thought to automatically activate a set of analog numerosity representations (Nieder et al., 2002; Nieder & Miller, 2003), as indicated by well-known distance and magnitude effects (Dehaene, Dupoux, & Mehler, 1990; Moyer & Landauer, 1967) and numerical Stroop effects (Henik & Tzelgov, 1982). Studies of numerical averaging have shown that human observers also have a remarkable ability to identify and average symbolic numbers even under stringent processing constraints. For example, Brezis et al. (2015, Exp. 3) presented observers with a sequence of four to 16 two-digit numbers at a rate of ten

✉ David Rosenbaum
davidros28@gmail.com

¹ Tel Aviv University, Tel Aviv-Yafo, Israel

² Paris School of Economics and CNRS, Centre d’Economie de la Sorbonne, Paris, France

66 items/s, and asked participants to indicate the average on a
67 continuous scale. The results show that the estimation preci-
68 sion (the RMSD) improved with the length of the sequence,
69 indicating that observers did not use only a limited sample of
70 the sequence. These results were accounted for by a
71 population-pooling mechanism (in which encoding noise
72 would average out over items).

73 In a few other studies it has been shown that observers can
74 extract numerical information from arrays of numerical symbols
Q5 75 presented simultaneously. For example, Corbett et al. (2006)
76 have shown that observers are able to discriminate between
77 two circular arrays of six digits (comprising 2s and 5s), presented
78 simultaneously (for as short as 80 ms), the one that had a higher
79 average (more 5s) with an accuracy exceeding 80% (Exp. 1).
80 Critically, this discrimination was faster and more accurate with
81 arrays made of the 2 and 5 symbols than with p and q symbols,
82 and this speedup only took place when the numerical meaning
83 could be used as a basis for the classification task (Exp. 3).
84 Finally, using a dual-task methodology, the authors have shown
85 that this ability requires central attention (Exp. 5). This study thus
86 demonstrates that numerical information is rapidly extracted
87 from arrays of numbers, at least when these arrays are relatively
88 simple. However, this very specific set of stimuli makes it possi-
89 ble for observers to adopt a strategy that might not involve
90 computing of an average over all elements.¹ Situations involving
91 more complex arrays remain to be investigated, as they might
92 help uncover the computational algorithms used by observers to
93 evaluate an average over items.

94 Two recent studies have taken this approach, using larger
95 arrays of two-digit numbers presented simultaneously (for up to
96 4–5 s) and asking participants to decide whether the average
97 was smaller or higher than a reference (Vandormael et al.,
98 2017; Vanunu et al., 2020). In these studies, the observers'
99 accuracy improved with presentation time, with the distance
100 of the average from the reference, and with sets involving lower
101 variance. The two studies differed, however, in their conclusion
102 about the algorithm used to carry out the task: whereas
103 Vandormael et al. (2017) found robust-averaging – an algo-
104 rithm that gives less weight to outliers, Vanunu et al. (2020)
105 found on the contrary that extreme values received equal or
106 higher weights. Although the reasons for this discrepancy are
107 still unclear, at least in both cases participants relied on some
108 items more than on others. This finding relates to the notion of
109 capacity that has been put forward in early cognitive models of
110 attention and working memory, and that has also been part of
111 recent theoretical accounts of ensemble perception.

112 In the context of extracting a set-average, capacity can be
113 defined as the number of items pooled together in the estimation

(Alik et al., 2013; Dakin, 2001; Solomon, May, & Tyler, 2016). 114Q6
Whereas this definition assumes an all-or-none selection of some 115
items and not others, one alternative view involving distributed 116
attention can be considered. In this view, all the elements con- 117
tribute to the estimation of the average, each element receiving a 118
fraction of the attentional resources available, which becomes 119
smaller when there are more elements in the array (Eriksen & 120Q7
StJames; Baek and Chong, 2020; Chong & Treisman, 2005). As 121 Q8
shown by Baek and Chong (2020a), a signature of this model is 122
an improved precision with set size (see also Brezis et al., 2015, 123
for the case of sequential presentation). 124

125 The appeal of the notion of capacity or distributed attention-
126 al resources is that these notions are domain general, and can be
127 compared across observers and across tasks. Surprisingly,
128 however, and despite the fact that many studies have
129 demonstrated that observers form ensemble representations
130 over various dimensions, how these dimensions compare, for
131 example in terms of the capacity, is not clear. In a recent study,
132 Haberman, Brady, and Alvarez (2015) found that individual
133 differences in performances (mean absolute errors when iden-
134 tifying the average over a set) were correlated between two
135 low-level features such as orientation and color, but uncorrelat-
136 ed when comparing a low-level feature to a higher-level feature
137 such as facial expression. This suggests that ensemble repre-
138 sentations for different features might operate with different
139 levels of performance, although capacity or distributed atten-
140 tion was not specifically assessed in this study.

141 Here, we hypothesize that the capacity with which observers
142 build an average representation might depend on how much
143 attentional and visual working memory resources are involved
144 in extracting and manipulating the task-relevant feature. For in-
145 stance, we expect that limitations in distributed attention or visual
146 working memory capacity (Cowan, 2001; Luck & Vogel, 1997)
147 will affect the averaging of symbolic numbers (as suggested by
148Q9 Corbett et al., 2006) more than of simple visual properties like
149 orientation, which can be processed pre-attentively (Braun &
150 Sagi, 1991; Treisman & Gelade, 1980) and which engage group-
151 ing and the formation of a holistic Gestalt (Hess & Field, 1999;
152 Kovács & Julesz, 1993). While there is some debate on the
153 capacity with which orientation can be averaged in a brief array
154 (Baek & Chong, 2020a; Dakin, 2001; Robitaille & Harris, 2011;
155 Solomon, May, & Tyler, 2016; see review in Baek & Chong,
156 2020b), we expect that capacity would be reduced for numerical
157 stimuli, which are likely to require more attentional resources
158 due to their higher visual complexity.

159 The aim of our study was to contrast averaging of numer-
160 ical and visual oriented elements, within the same observers,
161 and using an identical experimental design (with the same
162 visual display and response procedure for these two dimen-
163 sions). By manipulating the size of the item-set across trials,
164 we aimed to evaluate how performance changes with set size,
165 and reveal the capacity of the integration process. For both
166 dimensions (numbers vs. orientations) we asked participants

¹ One alternative is that the observers estimate if there are more 2s than 5s, but not by how much (which would allow to decide that the average is higher or lower than 3.5, but not by how much). Alternatively, observers may estimate the average, but this could be based on a VWM capacity sample of about four items.

167 to report an estimation of the average on a continuous scale, in
 168 order to encourage the integration of all items, and minimize
 169 the use of non-averaging heuristics that might arise in tasks
 170 based on a comparison to a reference. We expect that in the
 171 numerical averaging task, participants will be more accurate
 172 with smaller arrays. By contrast, in the orientation averaging
 173 task we expect either a fixed (or improved) precision with the
 174 set size of the array, as a result of averaging the encoding
 175 noise. To validate these conclusions, we used computational
 176 modeling to fit the data with two models, namely (1) the
 177 limited-capacity (subsampling) model (Alik et al., 2013;
 178 Solomon, May, & Tyler, 2016) and (2) the distributed atten-
 179 tion or ‘zoom lens’ model (Baek & Chong, 2020a), and ex-
 180 tracted the capacity or attention parameters for the two tasks.
 181 Finally, we examined the weights given to the mid-range and
 182 extreme values and compared them across the tasks
 183 (Vandormael et al., 2017; Vanunu et al., 2020).

184 **Experiment**

185 The experiment briefly presented arrays of numbers (digits
 186 1–9) or oriented elements (Gabors) of various set sizes (from
 187 four to 12) and required participants to estimate the numerical
 188 or orientation average on a continuous scale. We used an
 189 estimation on a continuous scale rather than a binary decision
 190 relative to a reference, as this minimizes the reliance on some
 191 heuristics, such as counting the number of elements higher
 192 than the reference, or even the number of extreme (high- vs.
 193 low-value) elements. Our main focus is the dependency of the
 194 estimation precision on set size in the two tasks.

195 **Methods**

196 **Participants**

197 Eighteen healthy adult volunteers with normal or corrected-to-
 198 normal vision participated in this study. All volunteers gave
 199 written informed consent to participate in this study. All pro-
 200 cedures and experimental protocols were approved by the
 201 ethics committee of the Psychology Department of Tel Aviv
 202 University (Application 743/12). All experiments were carried
 203 out in accordance with the approved guidelines. Due to the
 204 COVID-situation, testing conditions were restricted. We of-
 205 fered our participants the option to be tested (for an equivalent
 206 of \$15) in the lab under special safety COVID19 guidelines,
 207 or to run the experiment at home (same pay) from their own
 208 computer (to do this they needed to have Matlab installed on
 209 their computer). Ten participants were tested in the lab and
 210 eight were tested at home.

Stimuli

In the lab, displays were generated by an Intel I7 personal com-
 212 puter attached to a 24-in. Asus 248qe monitor with a 144-Hz
 213 refresh rate, using 1,920 × 1,080 resolution graphics mode. Due
 214 to the Covid19 situation, eight participants were tested at home
 215 using their own personal computers, but the experimental code
 216 was designed so as to detect the monitor’s resolution and present
 217 the stimuli with the same relative size. All participants were
 218 approximately at a distance of 60 cm from the screen.

The stimulus was an array of four, eight, or 12 elements
 220 (Gabor patches or numbers, depending on the task), randomly
 221 located on a gray background, within an invisible 5 × 6 grid
 222 (each cell was 77 × 96 pixels), with a restriction of no two
 223 horizontally adjacent elements and no element in the cells just
 224 above and below fixation (see Fig. 1a and b). Numbers were
 225 integers between 1 and 9, presented in white in David font size
 226 25. Gabor patches were 200 pixels wide, with a spatial frequen-
 227 cy of 0.2 cycles per pixel and standard deviation of 20 pixels.
 228 Gabors' orientations varied from 42° to 138° in nine equidistant
 229 steps. Stimuli were generated using Psychtoolbox for Matlab.
 230

Trial procedure

Each trial began with the onset of a central fixation dot (1 s)
 232 followed by the stimulus array (numbers or Gabors), which
 233 remained on the screen for 300 ms. After the offset of the
 234 array, participants were instructed to report the numbers' av-
 235 erage (number task) or the Gabors' average angle (orientation
 236 task) on a semicircular scale (an arc from 30° to 150°), using
 237 their mouse. The mouse cursor was always at the middle of the
 238 circle (red fixation dot) at the beginning of the scale display
 239 and thus it had an equal distance from each point of the scale
 240 (see Fig. 1). The scale labels were numbers from 1 to 9 (num-
 241 ber task) or oriented lines from 30° to 150° (orientation task).
 242 The participants had a 5-s deadline to respond.
 243

Design

Each participant completed both the orientation and the num-
 245 ber averaging tasks, in separate blocks of 360 trials each, in an
 246 order counterbalanced across participants. Set sizes (four,
 247 eight, 12 elements) were randomly interleaved across trials.
 248 Numbers were drawn randomly from one of three Gaussian
 249 distributions with means of 3.5, 5, or 6.5 and standard devia-
 250 tion of 1.5. The Gabors’ orientations were drawn from
 251 Gaussian distributions with means of 72°, 90°, and 108° rela-
 252 tive to horizontal, and SD of 18°. The positions of numbers
 253 1–9 on the response scale corresponded to orientations from
 254 42° to 138°. Due to a coding error, in four (of the 18) partic-
 255 ipants the Gaussian distribution of the Gabors were located at
 256 76.6°, 96.6°, and 116.6°, generating a small tilt of the overall
 257 distribution. This coding error was corrected in the other
 258

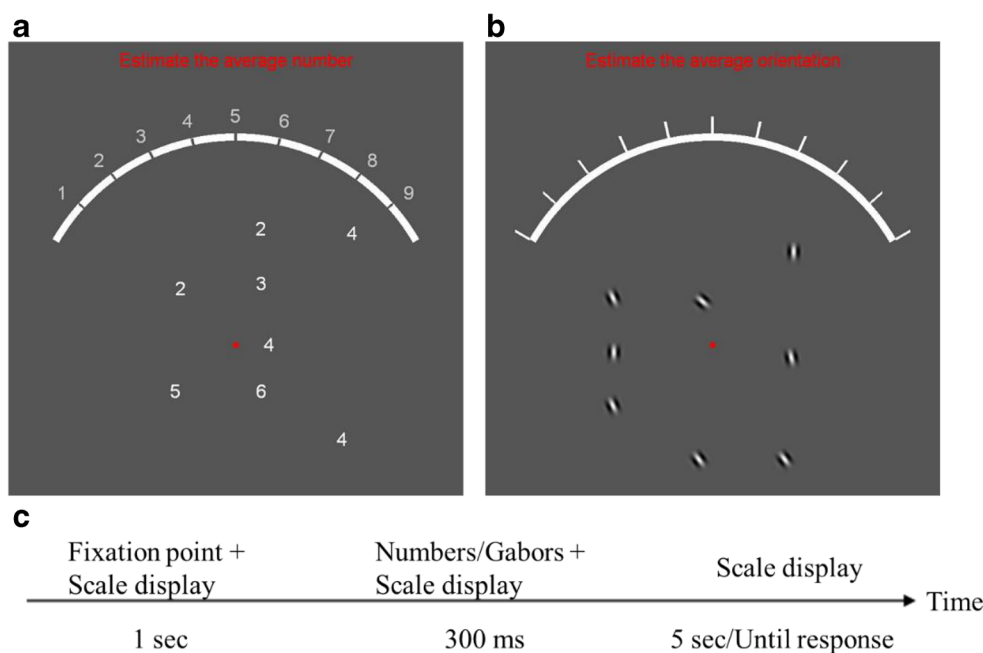


Fig. 1 Representative trial stimuli of each condition (set size 8). **(a)** Numbers condition. **(b)** Gabors condition. **(c)** Timeline diagram of a single trial. Each trial began with a fixation point for 1 s, followed by

259 participants. Since responses are made on a continuous scale
 260 and the actual deviation can be correctly extracted, all partic-
 261 ipants were included in the analysis.

262 Results

263 For simplicity and normalization between the two tasks we
 264 computed the different accuracy measurements in the orienta-
 265 tion task after we transformed the orientation angles to num-
 266 bers of 1–9, based on the mappings above.

267 Averaging precision

268 We used two measures to quantify participants' precision in
 269 the averaging tasks: First, we looked at the Pearson correlation
 270 across trials between the real and estimated averages of the
 271 array in each trial (see Fig. 2 for an example participant). The
 272 average correlation was high both for the orientation task (av-
 273 erage $r = .72$, $SD = 0.1$) and for the number task (average $r =$
 274 $.80$, $SD = 0.08$). Note that in both tasks we observed regres-
 275 sion to the mean, by which responses were biased towards the
 276 center of the scale. Second, we computed the root mean square
 277 deviation (RMSD) between the real averages and the partici-
 278 pants' responses across trials (see Fig. 3a). To obtain a chance-
 279 level baseline for this measure, we evaluated the RMSD for
 280 randomly shuffled responses across trials, both for the orienta-
 281 tion and the number tasks. We found the actual RMSD was
 282 significantly lower (more precise) than the shuffled version
 283 (orientation task: actual RMSD = 1.00, shuffled RMSD =

the array and ended with the response scale display. Trials end when the
 participant enters a response or after a 5-s deadline

1.84, $t(17) = 16.6$, $p < .001$. number task: actual RMSD = 284
 0.86, shuffled RMSD = 1.87, $t(17) = 25.9$, $p < .001$). 285

In order to test the main effect of set size and its interaction 286
 with task, we carried out a two-way repeated-measures 287
 ANOVA (set size \times task) with RMSD as the dependent vari- 288
 able. There was a significant interaction between the effects of 289
 set size and task, $F(2,34) = 6.5$, $p < .01$. A separate ANOVA for 290
 each task revealed a significant set size effect for the number 291
 task, $F(2,34) = 13.8$, $p < .001$, but not for the orientation task; 292
 $F(2,34) = 0.88$, $p = .68$. Post hoc comparison using Holm's test 293
 in the number task showed that RMSD was significantly lower 294
 (more precise) for four items than for eight and 12 items. In 295
 sum, in the number task participants were less accurate as set 296
 size increased, as opposed to the orientation task in which set 297
 size did not influence precision (see Fig. 3a). 298

299 Reaction times

To evaluate whether the decrease in performance with set size 300
 for numbers might be related to a potential speed-accuracy 301
 tradeoff, we also looked at response times (Fig. 3b). We re- 302
 peated the same two-way repeated-measures ANOVA (set 303
 size \times task) now with median response time (RT) as the de- 304
 pendent variable. There was a significant interaction between 305
 the effects of set size and task on RT, $F(2,34) = 7.7$, $p < .01$. A 306
 separate ANOVA for each task revealed a significant set size 307
 effect for the number task, $F(2,34) = 10.25$, $p < .001$, but not 308
 for the orientation task, $F(2,34) = 2.09$, $p = .13$. Post hoc com- 309
 parison using Holm's test in the number task showed that 310

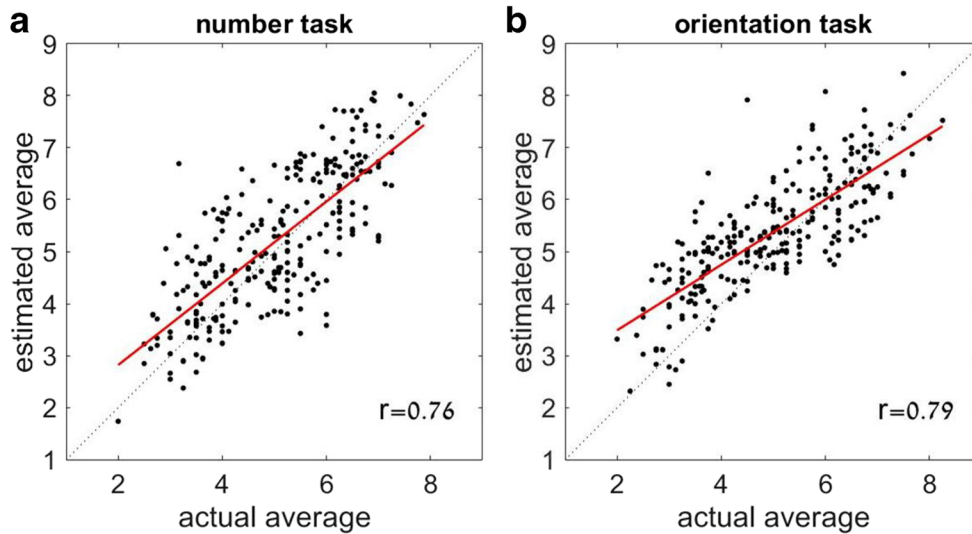


Fig. 2 (a) Correlation between the real average and the estimated average of a representative participant in the number task. (b) Correlation between the real average and the estimated average of a representative participant

in the orientation task. In both panels, each dot corresponds to a single trial, and the red line represents the regression of the estimated average against the actual average across trials

311 four-items RT was significantly slower than eight- and 12-
312 items RT (see Fig. 3b).

313 In order to understand if the slowdown in the number-
314 averaging task at the set size of four can account for the improved
315 precision in this condition, we computed for every participant the
316 correlation between absolute errors (RMSD) and RTs across
317 trials, separately for each set size and task. We reasoned that if
318 such a speed-accuracy tradeoff occurred, then longer RTs would
319 be associated with lower errors, resulting in negative correlations
320 between errors and RTs. However, no negative correlations were
321 found at the group level (see Fig. 1 in the Appendix for the
322 distribution of correlation coefficients across participants). In par-
323 ticular, for the numerical averaging task, the correlations were
324 close to zero (for set sizes four, eight, and 12 the mean *r* values

were -.012, -.004, and .042, respectively, with SDs 0.17, 0.14, 325
and 0.10, across participants). For the orientation averaging task 326
we found small (but statistically significant) positive correlations 327
at set sizes four (mean *r* = .092, SD = 0.15, *t*(17) = 2.59, *p* = .019) 328
and set size 12 (mean *r* = .057, SD = 0.084, *t*(17) = 2.88, *p* = .010). 329
To further discard the possibility of a speed-accuracy tradeoff for 330
the set size of four in the number task, we eliminated the 20% 331
slowest trials in that condition, so that the remaining trials had a 332
median RT that was the same as the set size eight condition, and 333
we examined RMSD in this RT-equivalent dataset. As expected 334
from the null correlation between RT and absolute errors, this 335
exclusion of slow trials did not affect the results regarding 336
RMSD. Critically, the interaction between set size and task was 337
maintained ($F(2, 34) = 6.13, p = .011$). 338

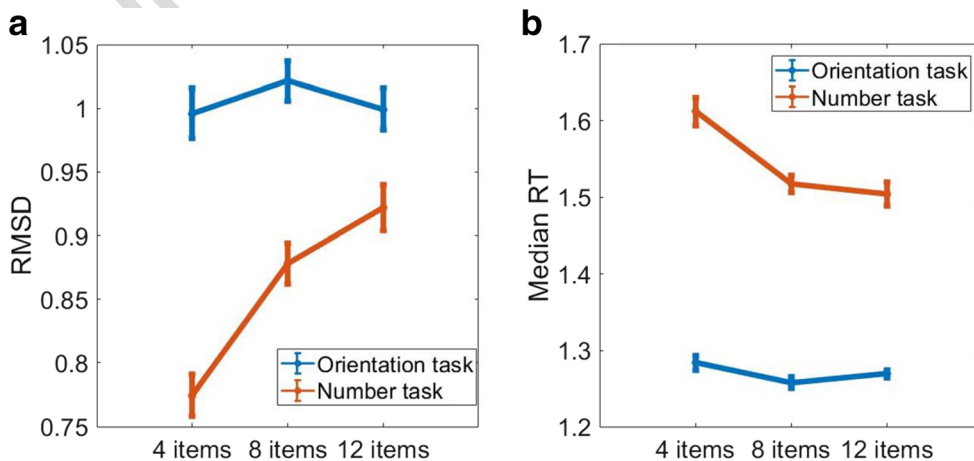


Fig. 3 (a) Root mean square deviation as a function of set size. In the orientation condition (blue) participants were not impacted by set size. In contrast, the number condition (red) shows that participants' performance deteriorated as set size increased. (b) Median response time (RT) as a function of set size. In the orientation task (blue) there was no difference

in RT between the different set sizes. In the number task (red) responses were slower in the four-items condition compared to the eight- and 12-items conditions. In both panels, errors bars represent the mean and its standard error across participants

339 **Summary and discussion**

340 While the participants were able to carry out both tasks relatively
 341 well (as indicated by correlations between real and estimated
 342 values higher than .70), the precision of their estimation showed
 343 a different dependency on the set size of the array in the two
 344 tasks. For orientation-averaging, set size did not affect either the
 345 precision or the mean RT, suggesting a parallel process (Ariely,
 346 2001; Chong & Treisman, 2005; Robitaille & Harris, 2011). For
 347 the numerical averaging on the other hand, both the precision
 348 and the RT decreased with set size. One possibility is that for
 349 small arrays (four digits), participants could have attempted to
 350 carry out the estimation by using a symbolic computation strat-
 351 egy, a strategy that they gave up on with larger arrays (Brezis
 352 et al., 2015). The null correlation between RMSD and RT ob-
 353 served in this condition indicates that this extra time did not help
 354 the participants to improve their estimation precision.

355 To conclude, we see that the ability of the participants to
 356 average larger arrays of numbers appears more limited, as the
 357 precision of the estimation is reduced with the size of the array.
 358 This is what would be expected if capacity (i.e., the number of
 359 elements the subjects can pool from) was reduced in the numer-
 360 ical task. In the next section we apply computational modeling
 361 in order to extract the capacity and attention parameters of the
 362 two tasks, and to examine additional biases, such as the weight
 363 given to in- or outlying elements (de Gardelle & Summerfield,
 364 2011; Vandormael et al., 2017; Vanunu et al., 2019).

365 **Computational analysis**

366 We applied two computational models to account for the data
 367 across all trials and participants, in both tasks. The first model
 368 is a version of the limited-capacity (subsampling) model (Alik
 369 et al., 2013; Dakin et al., 2001; Solomon, May, & Tyler,
 370 2016). This model assumes that out of N items presented, only
 371 M items are pooled up to generate the average-estimate. There
 372 are three sources of noise in this estimate. The first one is the
 373 sampling noise caused by subsampling (M out of N) elements.
 374 The second is an encoding noise, which is averaged out with
 375 M. The last component is a late-noise (this may include a
 376 motor component), which is not affected by M or N.

$$MeanEstimated = a + b \left(\frac{\sum_{i=1}^M x_i + \varepsilon_e}{M} \right) + \varepsilon_m, \varepsilon_e \sim N(0, \sigma_e^2) \text{ and } \varepsilon_m \sim N(0, \sigma_m^2) \quad (1)$$

379 The model is summarized by Eq. 1, where M is the number
 380 of sampled items out of the array, x_i is the ith item that was
 381 sampled, ε_e is the encoding noise, and ε_m is the motor noise. In
 382 this equation, a and b correspond to the intercept and slope
 383 parameters by which the internal estimation is mapped onto
 384 the external response-scale. Note that b < 1 would induce a

386 regression to the mean, which appears in the data (Fig. 2),
 387 and which is adaptive when observers face uncertainty but have
 388 prior knowledge about the distribution of the stimuli (Jazayeri
 389 & Shadlen, 2010; Anobile, Cicchini, & Burr, 2012).²

390 The second model is a version of the zoom lens model
 391 (Baek & Chong, 2020). The model assumes that while all
 392 visual elements contribute to the averaging estimation, they
 393 are subject to distributed attentional resources, which can vary
 394 from a sharp focus (for small arrays) to a broad one (for larger
 395 arrays). The precision of the processing is then in inverse
 396 proportion to the area of focus, similar to the zoom lens of a
 397 camera. As a result, the model assumes that an increase in set
 398 size leads to an increase in encoding noise for each item. There
 399 are also three sources of noise in this model. The first two are
 400 encoding noise and late noise, similar to the previous model.
 401 The third one is the attention parameter (A), which is a noise-
 402 reduction factor multiplied to encoding noise.

$$MeanEstimated = a + b \left(\frac{\sum_{i=1}^n x_i + \varepsilon_e}{n} \right) + \varepsilon_m, \quad (2)$$

$$\varepsilon_e \sim N(0, \sigma_e^2), \quad \varepsilon_m \sim N(0, \sigma_m^2) s$$

$$\sigma = \sqrt{\frac{(n-1+A)^2}{n^3} \sigma_e^2 + \sigma_m^2} \quad (3)$$

403 This model is summarized by Eqs. 2 and 3, where x_i is the
 404 ith item that in the array, n is the set size of the array, ε_e is the
 405 encoding noise, ε_m is the motor noise, A is the attention param-
 406 eter, and a and b correspond to the intercept and slope
 407 parameters by which the internal estimation is mapped onto
 408 the external response-scale.

409 Since fitting five parameters is computationally challeng-
 410 ing (from a model recovery perspective), we carried out the
 411 model fits in two steps. First, we conducted a simple regres-
 412 sion predicting the trial-by-trial response of each participant
 413 from the sequence-average, to determine the a and b param-
 414 eters for each participant. We then fixed those parameters and
 415 we fitted the three noise parameters, M, ε_e, ε_m, or A, ε_e, ε_m. For
 416 the zoom lens model, the predicted distribution of the estimat-
 417 ed mean is Gaussian around the actual average and with a
 418 variance determined by Eq. 3. For the sampling model, we
 419 resorted to simulations. For each trial we computed the ex-
 420 pected distribution of the estimated mean over the array, given
 421 the parameters of the model. In both models, from the predict-
 422 ed distribution (in each trial) we obtained the log-likelihood of
 423 the response of the observer in that trial. These log-likelihoods
 424 were accumulated across trials and the model parameters were
 425 optimized to maximize the total log-likelihood (see Tables 1
 426 and 2 in the Appendix for parameters AIC/BIC) Fig. 4.

² We also carried out model fitting without the a, b parameters, but the results were less good in terms of AIC/BIC measures, and they provide similar conclusions. So we only report the model comparisons that include intercept and slopes parameters.

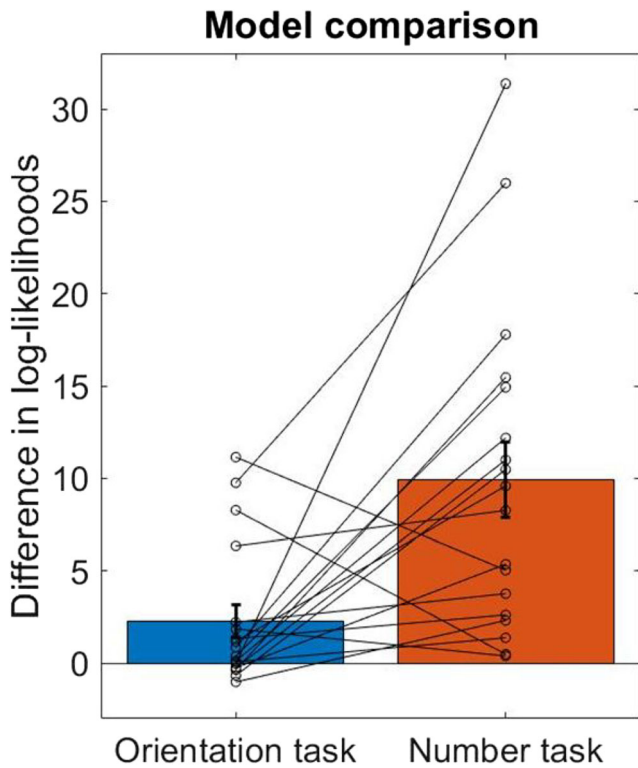


Fig. 4 The difference in log-likelihood between the zoom lens model and the sampling model as a function of task (orientation vs. numbers). Positive values indicate an advantage for the sampling model. Each dot is an individual observer. Error bars correspond to SEM

Model comparison

We compared the capacity/sampling and the zoom lens model to test which of them accounts better to the data in each task. The models have the same number of parameters so we compared directly the log-likelihoods. Figures 5 and 6 shows the difference in log-likelihood (zoom lens minus sampling model, such that positive values are in favor of the sampling model) in each task. As shown in the figure, there are very small

differences in the model fits in the orientation task (except for four subjects out of 18), but there are large differences in the number task, where the sampling model fares significantly better (see Tables 1 and 2 in the Appendix for more details).

Interestingly, all the participants for whom the sampling model wins over the zoom lens model are those for whom the fitted value of the capacity parameter, k , was very small (2 or 3; see Tables 1 and 2). Besides, all the participants for whom the capacity parameter was $k = 12$ in the orientation task (maximum value), were those for whom the zoom lens model won (see Tables 1 and 2). We next focused on how the capacity parameters vary with the task (see Figs. 5 and 6; see Tables 1 and 2 in the Appendix for other parameters). Despite marked variability across individuals, we observed overall a higher capacity in the orientation task ($M = 7.3, SD = 3.9$) than in the number task ($M = 3.6, SD = 1.8$). The difference between the two tasks was statistically significant ($t(17) = 3.5, p < .005$). This result confirms our hypothesis that when constructing their representation of the average over a set of items, observers integrate more items in the orientation task than in the number task.

Weights of inlying versus outlying elements

Finally, we examined the weights that participants gave to the different elements in the array, depending on their relative rank (among all elements in the array) and depending on the task (i.e., number or orientation). In particular, we compared elements falling in the middle of the sample (hereafter inlying elements) versus elements at the extreme (hereafter outlying elements). For example, for a sequence such as (2, 3, 4, 5, 5, 6, 7, 8), we considered that (2, 3, 7, 8) were outlying elements and that (4, 5, 5, 6) were inlying elements. For each task and set size, we then extracted the weights given to outlying and to inlying elements using the following linear regression (Eq. 2):

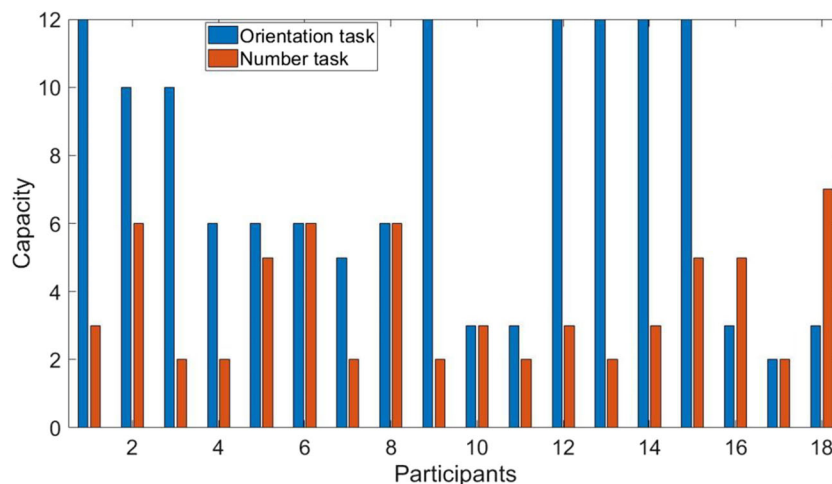


Fig. 5 Capacity-parameter M , for each of the participants in the two tasks

$$Response = \beta_0 + \beta_{in} \left(\frac{2}{n} \sum_{i \in In} X_i \right) + \beta_{out} \left(\frac{2}{n} \sum_{i \in Out} X_i \right)$$

476

477 with X_i the ordered samples, and $In = \llbracket \frac{n}{4} + 1, \frac{3n}{4} \rrbracket$ and Out
 478 $= \llbracket 1, \frac{n}{4} \rrbracket \cup \llbracket \frac{3n}{4} + 1, n \rrbracket$ the indices for inlying and outlying ele-
 479 ments, respectively.

480 We then examined how these weights varied across condi-
 481 tions (Figs. 5 and 6). A $2 \times 3 \times 2$ ANOVA (task, set size, in/
 482 outliers) shows a significant triple interaction ($F(2, 34) = 4.48$,
 483 $p = .031$). We thus conducted separate ANOVAs for each task,
 484 to examine the effect of set size and element rank. In the
 485 number task, there was only a main effect of rank, $F(1, 17)$
 486 $= 11.20$, $p = .004$, in which participants gave more weights to
 487 the outlying elements, in a similar manner across all set sizes.
 488 By contrast, for the orientation task there was both a main
 489 effect of set size, $F(2, 34) = 23.82$, $p < .001$, and an interaction
 490 between set size and rank, $F(2, 34) = 6.03$, $p = .011$. Further
 491 examination of this interaction indicated that inlying elements
 492 were down-weighted relative to outlying elements only for the
 493 largest sets (size 12: rank effect: $F(1, 17) = 10.44$, $p = 0.005$)
 494 but not for smaller sets (sizes four and eight: both $p > .05$).

495 **General discussion**

496 We examined and compared the ability of observers to esti-
 497 mate the average number and the average orientation of ele-
 498 ments presented simultaneously for a brief (300 ms) duration.
 499 Our experimental procedure required observers to make a re-
 500 sponse on a continuous scale, rather than a binary decision,
 501 and our results indicated that in both tasks the observers were

able, despite the presence of a regression to the mean compo- 502
 nent, to make good estimations (see Fig. 2). 503

The critical difference between the two tasks was the impact 504
 of the set size on the precision with which the average was 505
 estimated. We expected that the perception of numerical sym- 506
 bols would depend more on attentional and visual working 507
 memory resources, compared with the perception of oriented 508
 elements, which can generate a more holistic (texture) process 509
 (Dakin, 2001; Chong & Treisman, 2005; Robitaille & Haris, 510
 2011). We thus expected to find a higher capacity in the pooling 511
 of orientations compared with the pooling of numbers. These 512
 predictions were confirmed at the group level, using estimates of 513
 capacity based on a sampling model of averaging. In addition, 514
 we also compared this model to the (distributed attention) zoom 515
 lens model of averaging, which instead of sampling involved 516
 distributed attention over all elements (Baek & Chong, 2020b). 517Q16
 While in the orientation task the zoom lens and the sampling 518
 models were about equal in their fit performance, in the numeri- 519
 cal task the sampling model provided a better fit. Consistent 520
 with this, the estimation precision decreased with set size only in 521
 the numerical task and the extracted capacity parameter M was 522
 lower for the numerical task (average $M = 3.7$), compared to the 523
 orientation task (average $M = 7.3$). 524

In addition to these group differences, we also observed a 525
 large heterogeneity in both tasks. While some participants 526
 showed maximal capacity in the sampling model (M values 527
 that approached the maximum set size of 12) and RMSD 528
 decreasing with set size (as a result of efficient pooling), others 529
 showed low capacity (values of $M = 2$) and RMSD increasing 530
 with set size. This type of heterogeneity was previously re- 531
 ported for the orientation averaging (Solomon, May, & Tyler, 532
 2016). One possibility discussed by Solomon et al. (2016) is 533
 that the efficiency may be a function of expertise with the task. 534

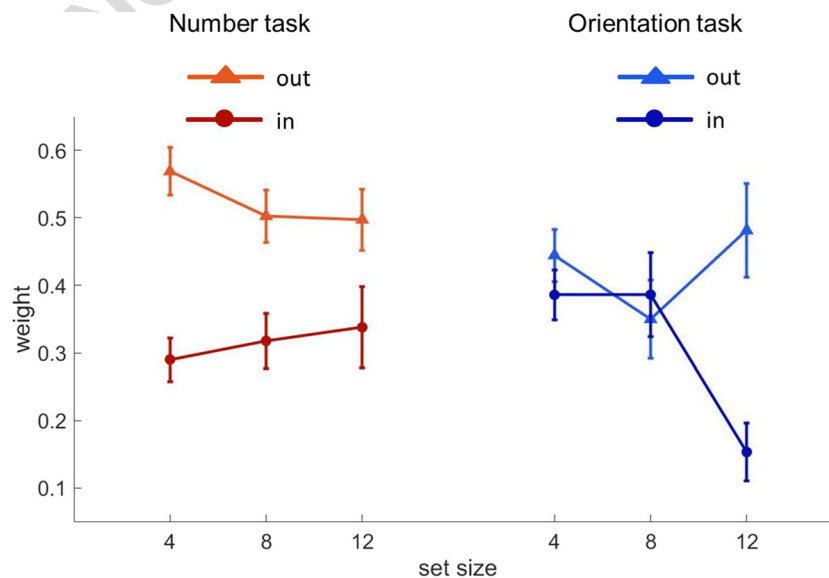


Fig. 6 Regression weights for inlying and outlying elements within each array, separately for the two tasks and the different set sizes. Error bars represents the mean and its standard error across participants

535 Our finding that orientation averaging is more efficient than
536 averaging of symbolic numbers is consistent with this possi-
537 bility: the visual system is arguably more expert in extracting
538 orientations from Gabor patches than in extracting the quan-
539 tity associated with a symbolic number. Could the inter-
540 individual variability in efficiency observed in our data also
541 relate to variations in expertise across participants?
542 Unfortunately, we cannot address this question directly with
543 our protocol, but there was room for variations in expertise
544 across participants, given that the amount of training our par-
545 ticipants received before engaging in the main experiment was
546 minimal (360 trials per task). In the case of averaging of sym-
547 bolic numbers in particular, one could further speculate that
548 familiarity with mathematics (e.g., due to studies, or to work-
549 related or other activities involving mental calculus) may af-
550 fect the efficiency with which participants compute an average
551 over a set of visually presented numbers. Future studies are
552 needed to further investigate this issue.

553 Our capacity estimate for the orientation averaging task is
554 somewhat higher than reported by Solomon et al. (2016) as
555 well as in some other studies (see, e.g., Table 1 in Solomon
556 et al., 2016). While as discussed above there was marked
557 heterogeneity in both studies, there are two aspects in the
558 experimental procedure that could account for potential dif-
559 ferences. First, while Solomon et al. (2016) used stimuli pre-
560 sented on a circular array, in our experiment they were pre-
561 sented in a texture type display, and random spatial positions,
562 which may enhance texture/grouping processes. Other studies
563 that used texture displays have also indicated a capacity that
564 exceeds the VWM of three to four items (Dakin, 2001;
565 Robitaille & Harris, 2011). Second, we used a continuous
566 response instead of a binary choice relative to a reference.
567 Doing this may have eliminated some non-integration strategy
568 to carry out the task, such as counting the elements higher than
569 the reference. Future work might investigate these aspects.

570 The main focus of our study was the comparison of the
571 capacity of the orientation and numerical averaging tasks.
572 Regarding this comparison, we should acknowledge one pot-
573 ential limitation of our experimental methodology, in that we
574 did not equate the visual characteristics of the stimuli between
575 the number task and the orientation task. It is possible that the
576 orientation stimuli may have benefited from a greater preci-
577 sion in terms of visual encoding than the number stimuli.
578 Indeed, our number stimuli involved higher spatial frequency
579 content (sharp edges), which may have been degraded to-
580 wards the periphery of the stimulus display. Fortunately, our
581 computational modeling allowed us to estimate encoding
582 noise for both the number task and the orientation task, and
583 it appears that irrespective of the model considered (sampling
584 vs. zoom lens), this early noise was actually higher for the
585 orientation task than in the number task (see Tables 1, 2, 3
586 and 4 in the Appendix Material), which we argue alleviates the
587 concern. Further research may, however, better address this

issue, by measuring the precision of the representation of single
items, in addition to the averaging task.

588
589
590 The lower capacity in the numerical averaging task indicates
591 that for most participants the estimation is based on sampling
592 only a few of the elements. Based on previous work (de Gardelle
593 & Summerfield, 2011; Vandormael et al., 2017; Vanunu et al.,
594 2020), we sought to investigate which elements received more
595 weight. The inlying/outlying analysis shown in Figs. 5 and 6
596 indicates that those elements are more likely to be extreme ele-
597 ments. Note that when a limited number of samples (say, two)
598 can be used for the averaging process, the precision of the esti-
599 mation is higher when the extreme ones are selected, compared
600 with a random selection. Thus, if these extreme elements are
601 easier to detect, relying on them could be an adaptive strategy.
602 This interpretation is consistent with the fact that in the orienta-
603 tion task, the weight of the extreme samples exceeds the weight
604 of the midrange samples, only at the largest set size (when the
605 set size exceeds the capacity of the orientation-averaging esti-
606 mation). While these results stand in contrast to those of
607 Vandormael et al. (2017), who reported robust averaging (lower
608 weights for extreme elements), they are consistent with those
609 reported by Vanunu et al. (2020). We should note that these two
610 studies used long presentation durations (several seconds in both
611 cases) and a binary comparison with a reference, whereas our
612 task involved brief displays and required an estimation on a
613 continuous scale.

614 The results for the numerical averaging also stand in con-
615 trast to those reported in Brezis et al. (2015, 2016, 2018), in
616 which the precision improved with set size, indicating pooling
617 across all (or almost all elements, from four to 16). The critical
618 difference, however, is that while in the present study, the
619 elements are briefly displayed simultaneously, in Brezis
620 et al. they were sequentially presented, resulting in less atten-
621 tional resource competition between the encoding of the ele-
622 ments. This suggests a framework in which while the estima-
623 tion mechanism is parallel (e.g., a neural population-coding
624 model in Brezis et al., 2016, 2018), the encoding of the items
625 has some serial (capacity limited) component that is lower for
626 symbols compared to oriented lines.

627 Finally, in addition to capacity, we also examined response
628 times for the two tasks. One interesting aspect was that RTs
629 markedly increased when participants had to average four
630 numbers, in comparison to eight or 12 numbers. Such an in-
631 crease for four elements was specific to the number task, and
632 did not occur in the orientation task. Thus, it might indicate
633 that participants approached the number averaging task differ-
634 ently with four items compared to eight or 12 items, for in-
635 stance by trying to calculate the average rather than by relying
636 on an intuitive estimation. We note, however, that these longer
637 response times did not lead to better responses. Whether this
638 change in strategy was deliberate or not and whether it may
639 reflect an adaptive strategy or not, however, remains to be
640 addressed. Future studies may investigate, in particular,

641 whether participants have a good insight or not about their
642 own cognitive processes in the averaging task.

643 **Appendix**
644 **Model fitting**

645 **Optimization procedure.** The free parameters of the sampling
646 and zoom lens models were fitted to the data of each participant
647 separately, using maximum likelihood estimation. We carried
648 out the model fits in two steps. First, we carried out a simple
649 regression predicting the trial by trial estimate of each subject
650 from the sequence-average, to determine the *a* and *b* parameters
651 for each subject. We then fixed those parameters and we con-
652 structed an *n*-dimensional grid (*n* is the number of free param-
653 eters for each model), with the four noise parameters (in total for
654 the two models), *M*, ϵ_e , ϵ_m , *A*. *M* ranging from 1 to 12 with
655 increments of 1, *A* ranging from 0 to 1 with increments of
656 .0101, ϵ_e ranging from 0 to 1.9 with increments of 0.126 for
657 the numbers task and ranging from 0 to 3 with increments of 0.2
658 for the orientation task and ϵ_m ranging from 0 to 1 with incre-
659 ments of 0.06 for the numbers task and ranging from 0 to 2 with
660 increments of 0.13 for the orientation task. This grid was
661 searched exhaustively, and for each set of parameters, θ_j , the
662 likelihood was calculated based on a Gaussian probability dis-
663 tribution function:

$$L(\theta_j) = \prod_{i=1}^N \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu_i)^2}{2\sigma^2}}$$

666 where *N* is the number of trials, x_i is the subject's estimated
667 average in each trial, μ_i is the predicted average by the model
668 excluding noise, and σ is the standard deviation such that σ
669 = $\sqrt{\sigma_e^2 + \sigma_m^2 + \sigma_M^2}$ for the efficiency model and
670 $\sigma = \sqrt{(n-1 + A) \frac{2}{n^3 \sigma_e^2 + \sigma_m^2}}$ for the Zoom lens model. We
671 also carried out model fitting without the *a*,*b* parameters and
672 compared the two fits.

673 **Model selection.** In order to evaluate the quantitative fits of
674 the models, we used two methods: (1) Akaike Information
675 Criterion (AIC; Akaike, 1973), and (2) Bayesian
Q17 676 Information Criterion (BIC; Schwarz, 1978; Raftery, 1995).
Q18 677 These selection criteria implement a trade-off between model
678 goodness of fit and complexity by penalizing additional free
679 parameters according to the following formulas:

$$AIC = -2 \cdot LL + 2 \cdot k$$

$$BIC = -2 \cdot LL + k \cdot \log(N)$$

683 where *LL* is the log-likelihood for the best fitting paramet-
684 ers, *k* is the number of free parameters and *N* is the number of
685 trials. AIC/BIC differences exceeding 10 are considered deci-
686 sive evidence in favor of the model with the lower numerical
687 values (Burnham & Anderson, 2002; Raftery 1995; see
688 Tables 1 and 2 for parameters and BIC/AIC values).

689 Fig. 7

Table 1 Comparison between the log-likelihood of the sampling model compared to the zoom lens model in the orientation task. The other columns show the parameters' value in each fit

Subject	Sampling model						Zoom lens model				t
	LL	a	b	M	ϵ_e	ϵ_m	LL	A	ϵ_e	ϵ_m	
1	422.9	1.19	0.84	12	2.2	0.80	422.8	0.31	3.0	0.6	t1.4
2	263.9	2.24	0.62	10	1.6	0.13	265.2	0.61	1.6	0.3	t1.5
3	442.6	1.28	0.75	6	2.6	0.40	442.7	0.30	2.2	1.0	t1.6
4	375.0	1.71	0.71	7	1.6	0.67	375.3	0.93	1.5	0.8	t1.7
5	350.4	1.2	0.78	6	1.2	0.27	352.6	0.00	0.7	0.6	t1.8
6	361.2	-0.03	0.98	6	1.4	0.00	361.3	0.00	1.6	0.4	t1.9
7	537.8	-0.02	1.1	5	2.2	0.00	538.6	0.00	1.8	0.9	t1.10
8	287.8	2.22	0.55	6	1	0.13	289.6	0.00	0.7	0.5	t1.11
9	271.5	2.74	0.46	12	0	0.53	271.1	0.00	0.3	0.5	t1.12
10	446.5	-0.3	1.07	3	1	0.13	456.3	0.00	0.9	0.8	t1.13
11	538.6	2.32	0.56	3	1.6	0.00	539.8	0.00	0.1	1.1	t1.14
12	563.4	-0.8	1.16	12	1.2	1.07	563.2	0.02	2.2	0.9	t1.15
13	438.5	1.82	0.62	12	1.2	0.67	438.2	0.01	2.0	0.5	t1.16
14	260.0	2.03	0.62	12	0.8	0.40	259.0	0.00	1.5	0.2	t1.17
15	390.6	3.23	0.38	12	1.8	0.27	389.9	0.16	2.2	0.0	t1.18
16	433.9	2.28	0.57	3	1	0.00	440.2	0.00	0.7	0.8	t1.19
17	544.2	1.57	0.65	2	1	0.00	555.4	0.00	0.5	1.1	t1.20
18	491.8	1.09	0.77	3	1.2	0.27	500.1	0.00	0.1	1.0	t1.21

Table 2 Comparison between the log-likelihood of the sampling model compared to the zoom lens model in the numbers task

Subject	Sampling model						Zoom lens model				t
	LL	A	b	M	ϵ_e	ϵ_m	LL	A	ϵ_e	ϵ_m	
1	239.7	0.3	0.9	3	0.1	0.2	271.0	0.00	0.1	0.7	t2.4
2	337.8	1.3	0.8	6	1.8	0.1	340.5	0.62	0.8	0.8	t2.5
3	332.5	-0.1	1.0	2	0.3	0.3	344.7	0.00	0.1	0.9	t2.6
4	325.3	0.6	0.9	2	0.1	0.3	340.3	0.00	0.8	0.8	t2.7
5	427.8	1.0	0.8	5	0.5	0.7	431.6	0.00	0.1	0.8	t2.8
6	343.2	1.1	0.8	6	1.3	0.0	344.5	0.07	0.6	0.6	t2.9
7	422.4	0.9	0.8	2	0.4	0.1	440.2	0.00	0.3	0.8	t2.10
8	361.6	1.3	0.7	6	1.3	0.2	362.0	0.21	0.8	0.6	t2.11
9	494.0	1.0	0.8	2	0.6	0.3	509.5	0.00	0.1	1.0	t2.12
10	276.0	0.3	1.0	3	0.1	0.1	301.9	0.00	0.1	0.6	t2.13
11	523.2	0.5	1.0	2	0.3	0.7	532.8	0.41	1.2	1.0	t2.14
12	443.6	0.6	0.9	3	0.8	0.4	448.9	0.00	0.6	0.8	t2.15
13	438.8	1.0	0.8	2	0.4	0.1	449.8	0.00	0.7	0.8	t2.16
14	305.5	1.1	0.8	3	0.4	0.1	307.8	0.00	0.8	0.5	t2.17
15	305.9	1.1	0.9	5	0.9	0.2	316.3	0.00	0.1	0.6	t2.18
16	399.0	0.9	0.8	5	0.5	0.6	407.3	0.00	0.7	0.7	t2.19
17	562.0	1.6	0.7	2	0.8	0.6	567.0	0.00	0.1	1.2	t2.20
18	475.6	1.7	0.6	7	1.1	0.7	476.1	1.00	0.3	0.9	t2.21

The other columns show the parameters' value in each fit

Atten Percept Psychophys

t3.1 **Table 3** Comparison between the AIC and BIC parameters for the sampling model with the mapping parameters compared to without the mapping parameters in the orientation task. The other columns show the parameters' value in each fit

t3.2	Subject	Model with mapping parameters						Model without mapping parameters					
t3.3		BIC	AIC	a	b	M	ϵ_e	ϵ_m	BIC	AIC	M	ϵ_e	ϵ_m
t3.4	1	875	856	1.19	0.84	12	2.2	0.80	902	890	12	2.6	0.80
t3.5	2	557	538	2.24	0.62	10	1.6	0.13	735	724	7	1.4	0.67
t3.6	3	915	895	1.28	0.75	6	2.6	0.40	919	907	2	1.2	0.53
t3.7	4	779	760	1.71	0.71	7	1.6	0.67	820	809	6	1	0.93
t3.8	5	730	711	1.2	0.78	6	1.2	0.27	784	773	5	1.2	0.27
t3.9	6	752	732	-0.03	0.98	6	1.4	0.00	750	739	6	1.4	0.00
t3.10	7	1,105	1,086	-0.02	1.1	5	2.2	0.00	1166	1155	2	1	0.53
t3.11	8	605	586	2.22	0.55	6	1	0.13	851	839	6	1.6	0.13
t3.12	9	572	553	2.74	0.46	12	0	0.53	926	914	6	0	0.80
t3.13	10	922	903	-0.3	1.07	3	1	0.13	920	908	3	0.8	0.40
t3.14	11	1,107	1,087	2.32	0.56	3	1.6	0.00	1175	1163	3	1.6	0.53
t3.15	12	1,156	1,137	-0.8	1.16	12	1.2	1.07	1,155	1,143	12	1.8	0.93
t3.16	13	906	887	1.82	0.62	12	1.2	0.67	1003	991	11	0.4	0.93
t3.17	14	549	530	2.03	0.62	12	0.8	0.40	796	784	6	1.2	0.40
t3.18	15	811	791	3.23	0.38	12	1.8	0.27	1,084	1,073	12	0	1.07
t3.19	16	897	878	2.28	0.57	3	1	0.00	1,028	1,017	2	0.4	0.53
t3.20	17	1,118	1,098	1.57	0.65	2	1	0.00	1,161	1,149	2	1	0.40
t3.21	18	1,013	994	1.09	0.77	3	1.2	0.27	1,038	1,026	3	1.2	0.40

t4.1 **Table 4** Comparison between the AIC and BIC parameters for the sampling model with the mapping parameters compared to without the mapping parameters in the numbers task

t4.2	Subject	Model with mapping parameters						Model without mapping parameters					
t4.3		BIC	AIC	a	b	M	ϵ_e	ϵ_m	BIC	AIC	M	ϵ_e	ϵ_m
t4.4	1	509	489	0.30	0.92	3	0.13	0.20	506	494	3	0.38	0.00
t4.5	2	705	686	1.25	0.78	6	1.77	0.13	730	718	5	1.65	0.33
t4.6	3	694	675	-0.11	1.03	2	0.25	0.27	684	672	2	0.25	0.27
t4.7	4	680	661	0.55	0.94	2	0.13	0.27	681	669	2	0.38	0.07
t4.8	5	885	866	0.95	0.80	5	0.51	0.67	911	899	5	1.52	0.27
t4.9	6	716	696	1.13	0.79	6	1.27	0.00	775	763	5	1.27	0.13
t4.10	7	874	855	0.94	0.83	2	0.38	0.07	868	857	2	0.13	0.00
t4.11	8	753	733	1.34	0.69	6	1.27	0.20	862	851	11	1.01	0.67
t4.12	9	1,018	998	1.03	0.82	2	0.63	0.27	1,031	1019	2	0.25	0.53
t4.13	10	581	562	0.31	0.95	3	0.13	0.07	572	560	3	0.00	0.00
t4.14	11	1,076	1,056	0.54	0.99	2	0.25	0.67	1,112	1100	1	0.00	0.07
t4.15	12	917	897	0.60	0.89	3	0.76	0.40	917	905	3	0.89	0.33
t4.16	13	907	888	1.02	0.84	2	0.38	0.13	921	909	2	0.00	0.27
t4.17	14	640	621	1.08	0.83	3	0.38	0.13	713	701	3	0.00	0.33
t4.18	15	641	622	1.06	0.85	5	0.89	0.20	751	740	3	0.25	0.33
t4.19	16	827	808	0.88	0.78	5	0.51	0.60	872	861	3	0.38	0.53
t4.20	17	1,153	1,134	1.62	0.65	2	0.76	0.60	1,192	1180	2	1.27	0.13
t4.21	18	981	961	1.72	0.63	7	1.14	0.73	1,053	1042	5	0.51	0.93

The other columns show the parameters' value in each fit

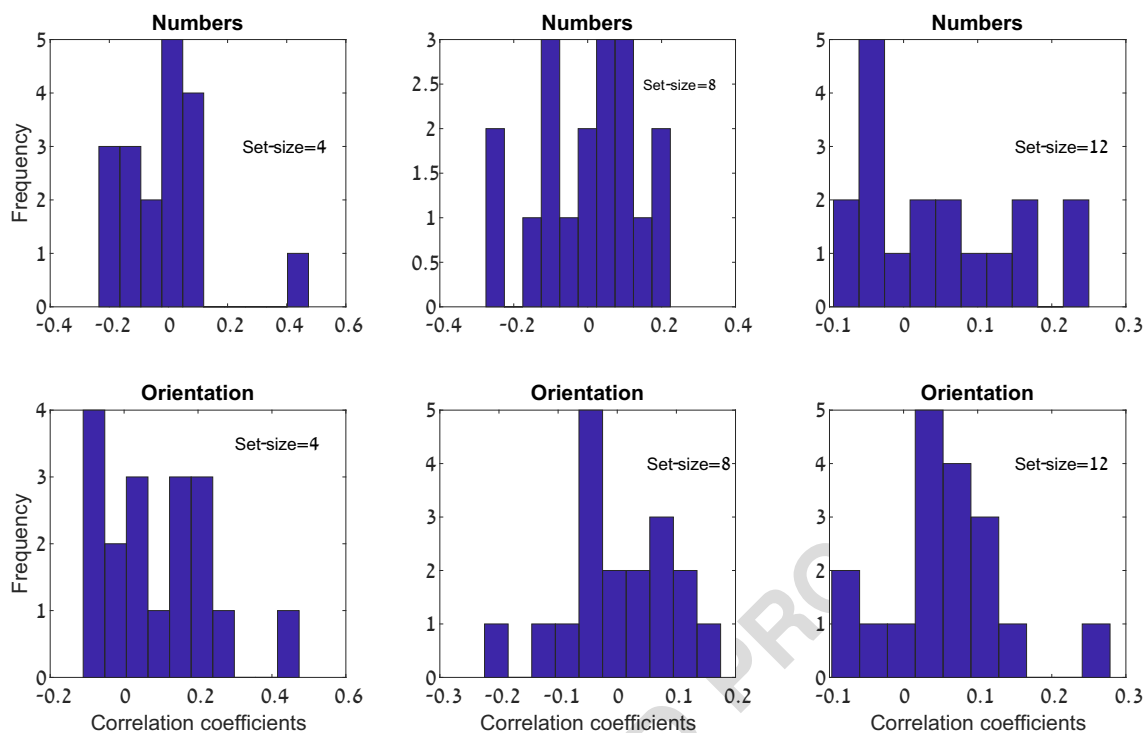


Fig. 7 Distribution of correlation coefficients across participants between root means square deviation (RMSD) and response time (RT)

690 **References**

691 Akaike, H. (1973). Information theory and an extension of the maximum
 692 likelihood principle. In BN, Petrov & F. Csaki (Eds.), Proceedings of
 693 the 2nd International Symposium on Information Theory (pp. 267-
 694 281). Budapest: Akademiai Kiado.
 695 Anobile, G., Cicchini, G. M., & Burr, D. C. (2012). Linear mapping of
 696 numbers onto space requires attention. *Cognition*, 122(3), 454-459.
 697 Ariely, D. (2001). Seeing Sets: Representation by Statistical Properties.
 698 *Psychological Science*, 12(2), 157-162. [https://doi.org/10.1111/
 699 1467-9280.00327](https://doi.org/10.1111/1467-9280.00327)
 700 Baek, J., & Chong, S. C. (2020). Distributed attention model of percep-
 701 tual averaging. *Attention, Perception, & Psychophysics*, 82(1), 63-
 702 79.
 703 Braun, J., & Sagi, D. (1991). Texture-based tasks are little affected by
 704 second tasks requiring peripheral or central attentive fixation.
 705 *Perception*, 20(4), 483-500.
 706 Brezis, N., Bronfman, Z. Z., Jacoby, N., Lavidor, M., & Usher, M.
 707 (2016). Transcranial direct current stimulation over the parietal cortex
 708 improves approximate numerical averaging. *Journal of
 709 Cognitive Neuroscience*, 28(11), 1700-1713.
 710 Brezis, N., Bronfman, Z. Z., & Usher, M. (2015). Adaptive Spontaneous
 711 Transitions between Two Mechanisms of Numerical Averaging.
 712 *Nature Publishing Group*, 1-11. <https://doi.org/10.1038/srep10415>
 713 Brezis, N., Bronfman, Z. Z., & Usher, M. (2018). A perceptual-like pop-
 714 ulation-coding mechanism of approximate numerical averaging.
 715 *Neural Computation*, 30(2), 428-446.
 716 Chong, S. C., & Treisman, A. (2003). Representation of statistical prop-
 717 erties. *Vision research*, 43(4), 393-404.
 718 Chong, S. C., & Treisman, A. (2005). Attentional spread in the statistical
 719 processing of visual displays. *Perception and Psychophysics*, 67(1),
 720 1-13. <https://doi.org/10.3758/BF03195009>
 721 Cowan, N. (2001). The magical number 4 in short-term memory: A
 722 reconsideration of mental storage capacity. *Behavioral and brain
 723 sciences*, 24(1), 87-114.

Dakin, S. C. (2001). Information limit on the spatial integration of local
 orientation signals. *JOSA A*, 18(5), 1016-1026. 724
 De Gardelle, V., & Summerfield, C. (2011). Robust averaging during
 perceptual judgment. *Proceedings of the National Academy of
 Sciences of the United States of America*, 108(32), 13341-13346. 725
 10.1073/pnas.1104517108 726
 Dehaene, S., Dupoux, E., & Mehler, J. (1990). Is Numerical
 Comparison Digital? Analogical and Symbolic Effects in Two-
 Digit Number Comparison. *Journal of Experimental Psychology*,
 16(3), 626-641. 727
 Haberman, J., Brady, T. F., & Alvarez, G. A. (2015). Individual differ-
 ences in ensemble perception reveal multiple, independent levels of
 ensemble representation. *Journal of Experimental Psychology:*
 General, 144(2), 432. 728
 Haberman, J., & Whitney, D. (2011). Efficient summary statistical rep-
 resentation when change localization fails. *Psychonomic Bulletin
 and Review*, 18(5), 855-859. [https://doi.org/10.3758/s13423-011-
 0125-6](https://doi.org/10.3758/s13423-011-0125-6) 729
 Hess, R. F., & Field, D. J. (1999). Integration of contours: New insights.
Trends in Cognitive Sciences, 3(12), 480-486. 730
 Jazayeri, M., & Shadlen, M. N. (2010). Temporal context calibrates inter-
 val timing. *Nature neuroscience*, 13(8), 1020. 731
 Khayati, N., & Hochstein, S. (2018). Perceiving set mean and range:
 Automaticity and precision. *Journal of Vision*, 18(9), 23-23 732
 Kovács, I., & Julesz, B. (1993). A closed curve is much more than an
 incomplete one: Effect of closure in figure-ground segmentation.
*Proceedings of the National Academy of Sciences of the United
 States of America*, 90(16), 7495-7497 733
 Luck, S. J., & Vogel, E. K. (1997). The capacity of visual working
 memory for features and conjunctions. *Nature*, 390(6657), 279-281. 734
 Moyer, R. S., & Landauer, T. K. (1967). Time required for judgements of
 numerical inequality. *Nature*, 215(5109), 1519-1520. 735
 Parkes, L., Lund, J., Angelucci, A., Solomon, J. A., & Morgan, M.
 (2001). Compulsory averaging of crowded orientation signals in
 human vision. *Nature neuroscience*, 4(7), 739-744. 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758

- 759 Robitaille, N., & Harris, I. M. (2011). When more is less: Extraction of
760 summary statistics benefits from larger sets. *Journal of Vision*,
761 11(12), 1–8. <https://doi.org/10.1167/11.12.18>
- 762 Sato, H., & Motoyoshi, I. (2020). Distinct strategies for estimating the
763 temporal average of numerical and perceptual information. *Vision*
764 *Research*, 174, 41–49.
- 765 Solomon, J. A., May, K. A., & Tyler, C. W. (2016). Inefficiency of
766 orientation averaging: Evidence for hybrid serial/parallel temporal
767 integration. *Journal of vision*, 16(1), 13–13.
- 768 Spitzer, B., Waschke, L., & Summerfield, C. (2017). Selective
769 overweighting of larger magnitudes during noisy numerical compar-
770 ison. *Nature Human Behaviour*, 1(8), 1–8.
- 771 Treisman, A. M., & Gelade, G. (1980). A feature-integration theory of
772 attention. *Cognitive psychology*, 12(1), 97–136.
- 784
- Vandormael, H., Herce, S., Balaguer, J., Li, V., & Summerfield, C. 773
(2017). Robust sampling of decision information during perceptual
774 choice. <https://doi.org/10.1073/pnas.1613950114> 775
- Van Opstal, F., & Verguts, T. (2011). The origins of the numerical dis- 776
tance effect: the same–different task. *Journal of Cognitive* 777
Psychology, 23(1), 112–120. 778
- Vanunu, Y., Hotaling, M., & Newell, R. (2020). Elucidating the differ- 779
ential impact of extreme-outcomes in perceptual and preferential 780
choice. *Cognitive Psychology*. 781
- Publisher's note** Springer Nature remains neutral with regard to jurisdic- 782
tional claims in published maps and institutional affiliations. 783

UNCORRECTED PROOF

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES.

- Q1. Please check if the author group is presented correctly.
- Q2. Ref. "Corbett, Oriet, & Rensink, 2006" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q3. Ref. "Nieder et al., 2002; Nieder & Miller, 2003" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q4. Ref. "Henik & Tzelgov, 1982" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q5. Ref. "Corbett et al. (2006)" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q6. Ref. "Alik et al., 2013" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q7. Ref. "Eriksen & StJames" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q8. The citation "Baek & Chong, 2020 a,b" has been changed to "Baek and Chong, 2020" to match the author name/date in the reference list. Please check if the change is fine in this occurrence and modify the subsequent occurrences, if necessary.
- Q9. Ref. "Corbett et al., 2006" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q10. The citation "Arieli, 2001" has been changed to "Ariely, 2001" to match the author name/date in the reference list. Please check if the change is fine in this occurrence and modify the subsequent occurrences, if necessary.
- Q11. Ref. "Vanunu et al., 2019" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q12. Ref. "Alik et al., 2013" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q13. Ref. "Dakin et al., 2001" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q14. Missing citation for Figure 4 was inserted here. Please check if appropriate. Otherwise, please provide citation for Figure 4. Note that the order of main citations of figures/tables in the text must be sequential.
- Q15. Figure5 here hasbeen changed to Figures 5 and 6 . Kindly check if appropriate.
- Q16. Ref. "Baek & Chong, 2020" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q17. Ref. "Schwarz, 1978" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.
- Q18. Ref. "Raftery, 1995" is cited in the body but its bibliographic information is missing. Kindly provide its bibliographic information in the list.