

---

*Research Article: New Research | Cognition and Behavior*

## **Differential involvement of EEG oscillatory components in sameness vs. spatial-relation visual reasoning tasks**

<https://doi.org/10.1523/ENEURO.0267-20.2020>

**Cite as:** eNeuro 2020; 10.1523/ENEURO.0267-20.2020

Received: 15 June 2020

Revised: 20 October 2020

Accepted: 21 October 2020

---

*This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.*

**Alerts:** Sign up at [www.eneuro.org/alerts](http://www.eneuro.org/alerts) to receive customized email alerts when the fully formatted version of this article is published.

Copyright © 2020 Alamia et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 **Differential involvement of EEG oscillatory components in sameness**  
2 **vs. spatial-relation visual reasoning tasks**

3 Andrea Alamia<sup>1</sup>, Canhuang Luo<sup>1</sup>, Matthew Ricci<sup>2</sup>, Junkyung Kim<sup>2</sup>, Thomas Serre<sup>2,3</sup>  
4 and Rufin VanRullen<sup>1,3</sup>

5 [rufin.vanrullen@cns.fr](mailto:rufin.vanrullen@cns.fr)

6  
7 <sup>1</sup>CerCo, CNRS Université de Toulouse, Toulouse 31055 (France)

8 <sup>2</sup>Department of Cognitive, Linguistic & Psychological Sciences; Carney Institute for Brain  
9 Science; Brown University, Providence, RI 02912, USA.

10 <sup>3</sup>ANITI, Université de Toulouse, Toulouse 31055 (France)  
11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26 **Abstract**

27 The development of deep convolutional neural networks (CNNs) has recently led to great  
28 successes in computer vision and CNNs have become de facto computational models of  
29 vision. However, a growing body of work suggests that they exhibit critical limitations beyond  
30 image categorization. Here, we study one such fundamental limitation, for judging whether  
31 two simultaneously presented items are the same or different (SD) compared to a baseline  
32 assessment of their spatial relationship (SR). In both human subjects and artificial neural  
33 networks, we test the prediction that SD tasks recruit additional cortical mechanisms which  
34 underlie critical aspects of visual cognition that are not explained by current computational  
35 models. We thus recorded EEG signals from human participants engaged in the same tasks  
36 as the computational models. Importantly, in humans the two tasks were matched in terms of  
37 difficulty by an adaptive psychometric procedure: yet, on top of a modulation of evoked  
38 potentials, our results revealed higher activity in the low beta (16-24Hz) band in the SD  
39 compared to the SR conditions. We surmise that these oscillations reflect the crucial  
40 involvement of additional mechanisms, such as working memory and attention, which are  
41 missing in current feed-forward CNNs.

42  
43 **Keywords:** Visual reasoning, spatial relationship, EEG oscillations, ERPs, deep neural  
44 networks.

45 **Significance statement**

46 Convolutional neural networks (CNNs) are currently the best computational models of  
47 primate vision. Here, we independently confirm prior results suggesting that CNNs can learn  
48 to solve visual reasoning problems involving spatial relations much more easily than  
49 problems involving sameness judgments. We hypothesize that these results reflect different  
50 computational demands between the two tasks and conducted a human EEG experiment to  
51 test this hypothesis. Our results suggest a significant difference – both in evoked potentials  
52 and in the oscillatory dynamics– of the EEG signals measured from human participants  
53 performing these two tasks. We interpret this difference as the signature for the fundamental  
54 involvement of recurrent mechanisms implementing cognitive functions such as working  
55 memory and attention.

56

57

58

## 59 1. Introduction

60 The field of artificial vision witnessed an impressive boost in the last few years, driven  
61 by the striking results of deep convolutional neural networks (CNNs). Such hierarchical  
62 neural networks process information sequentially – through a feedforward cascade of  
63 filtering, rectification and normalization operations. The accuracy of these architectures is  
64 now approaching – sometimes exceeding – that of human observers on key visual  
65 recognition tasks including object (He, Zhang, Ren, & Sun, 2016) and face recognition (P. J.  
66 Phillips et al., 2018). These advances suggest that purely feedforward mechanisms suffice to  
67 accomplish remarkable results in object categorization, in line with previous experimental  
68 studies on humans (VanRullen & Thorpe, 2001) and animals (Hollard & Delius, 1982;  
69 Vogels, 1999). However, despite the remarkable accuracy reached in these recognition  
70 tasks, the limitations of CNNs are becoming increasingly evident (see Serre, 2019 for a  
71 recent review). Beyond image categorization tasks, CNNs appear to struggle to learn to  
72 solve relatively simple visual reasoning tasks otherwise trivial for the human brain (Kim,  
73 Ricci, & Serre, 2018; Stabinger, Rodríguez-Sánchez, & Piater, 2016). A recent study (Kim et  
74 al., 2018) thoroughly investigated the ability of CNN architectures to learn to solve various  
75 visual reasoning tasks, and found an apparent dichotomy between two sorts of problems: on  
76 the one hand, tasks that require judging the spatial relations between items (Spatial  
77 Relationship – SR); on the other, those that require comparing items (Same-Different – SD).  
78 Importantly, Kim and colleagues demonstrated that CNNs can more easily learn the first  
79 class of problems compared to the second one.

80 This prompts the question of how biological visual systems handle such tasks so  
81 efficiently. Kim et al. (2018) suggest that SR and SD tasks tap into distinct computational  
82 mechanisms, thus leading to the prediction that different cortical processes are also involved  
83 when humans perform the two tasks: SR tasks can be successfully solved by feedforward  
84 processes, whereas SD tasks seem to require additional computations, such as working  
85 memory and attention. Here, we tested this hypothesis in two steps: first, we confirmed and  
86 extended Kim's results by comparing the performance of CNNs on an experiment in which  
87 we directly contrasted SD and SR tasks on the same stimulus set. Second, we recorded  
88 electrophysiological responses (EEG) in healthy human participants for the same  
89 experiment, after having matched the difficulty level via an adaptive psychometric procedure.  
90 We hypothesized that the additional computations required by the SD task, as compared to  
91 SR tasks, would elicit differences in evoked potentials (e.g. P300 modulations, which have  
92 been related to attentional mechanisms (Nash & Fernandez, 1996)) and brain rhythms  
93 related to working memory (such as beta-band oscillations (Benchenane, Tiesinga, &  
94 Battaglia, 2011; Lundqvist, Herman, Warden, Brincat, & Miller, 2018)). We found indeed that,

95 in addition to a variation in evoked potentials, the SD task elicited higher activity in specific  
96 beta-band oscillatory components in the occipital-parietal areas, which are typically  
97 associated with attention- and memory-related processes. We emphasize that the goal of the  
98 present study was not to identify the precise neural computations involved in the two tasks  
99 (which would naturally require a broader experimental set-up than a single EEG study), but  
100 rather to validate the hypothesis that SD involves additional computations relative to SR  
101 (even when the two tasks are equally difficult). We hope that this demonstration can be a first  
102 step towards characterizing the processes taking place in visual cortex during visual  
103 reasoning tasks, and designing more reliable and more human-like computational models.

## 104 **2. Materials & Methods**

### 105 **2.1 Participants and pilot experiment**

106 Twenty-eight participants (aged 21–34 years old with a mean age of  $26.6 \pm 3.7$ , 11  
107 women, 5 left-handed), volunteered to join the experiment. All subjects reported normal or  
108 corrected to normal vision and had no history of epileptic seizures or neurological disorders.  
109 Participants were pooled in two groups of 14 each: one group performed a pilot experiment,  
110 while the second one was tested on a final version of the task. The only difference between  
111 the pilot and the main study was the QUEST adaptive procedure used to match the difficulty  
112 level between conditions, which was not implemented in the pilot experiment. However, in  
113 both studies we found the very same result (see below, specifically fig. 4 and 5). In the main  
114 experiment, we kept the same number of participants to replicate the effect, after having  
115 removed the behavioral difference in task difficulty via the QUEST algorithm. This study  
116 complies with the guidelines of the research center where it was carried out, and the protocol  
117 was approved by an external committee (ethics approval number N° 2016-A01937-44). All  
118 participants gave written informed consent before starting the experiment, in accordance with  
119 the Declaration of Helsinki, and received monetary compensation for their participation.

### 120 **2.2 Experimental design**

121 The experiment was composed of 16 experimental blocks of 70 trials each, with a  
122 total duration of about 1 hour. Each trial lasted  $\sim 2$  seconds (Fig. 1A): 350ms after the onset  
123 of a black fixation cross ( $0.6^\circ$  width), 2 shapes were displayed for 30ms on opposite sides of  
124 the screen, distant  $2 \cdot \rho$  from each other with an angle of  $\pm(45^\circ + \theta)$  with respect to the  
125 horizontal midline ( $\rho$  being the distance from the center of the screen, and  $\theta$  the angular  
126 difference with the diagonal; see Fig. 1B). Each shape was selected from a subset of 36  
127 hexominoes, a geometric figure composed of 6 contiguous squares (see Fig. 1B) One  
128 second after the onset of the hexominoes, the fixation cross turned blue, cuing participants to  
129 respond. In half of the blocks, participants had to report whether the two shapes were the

130 same or different (Same-Different – SD condition); in the remaining blocks participants had to  
131 judge whether the two stimuli were aligned more horizontally or vertically (Spatial Relation –  
132 SR condition). Shapes were displayed at opposite sides of the screen along two main  
133 possible orientation axes sampled at random for every trial (either  $45^\circ$  and  $225^\circ$  or  $-45^\circ$  and -  
134  $225^\circ$ ). Both stimuli positions were jittered by a random offset  $\Delta x$  and  $\Delta y$  in both the x and y  
135 axis and a rotation  $\theta$  from the main axis. The same offsets were applied to both shapes, so  
136 they did not affect the angle between stimuli. The aim of such offsets was to prevent  
137 participants in the SR condition from determining the configuration of the two stimuli  
138 (orientation task) by merely judging the position of a single stimulus: without the random  
139 offsets, considering for example the top-right corner position, if the item were below/above  
140 the (imaginary) screen diagonal line, the overall orientation would be horizontal/vertical,  
141 without the need to consider the position of the corresponding bottom-left item. The offset  
142 then compelled participants to consider the relative position of both hexominoes at once.  
143 Importantly, in the main experiment (compared to the pilot experiment) the difficulty of the  
144 two tasks was controlled by an adaptive psychometric procedure (QUEST method, Watson &  
145 Pelli, 1983), which varied the eccentricity of the two stimuli  $\rho$  (in the SD blocks) or  $\theta$  (in the  
146 SR blocks) to maintain an overall accuracy level of 80% throughout the whole experiment. In  
147 fact, larger (smaller) values of  $\rho$  made the stimuli more (less) eccentric and the task more  
148 (less) difficult; similarly, smaller (larger) values of  $\theta$  set the stimuli closer to (farther from) the  
149  $45^\circ$  diagonal line, making the task more (less) difficult. We modified one parameter per  
150 condition (i.e., per block), while the other was kept constant (using the same value as in the  
151 preceding block). After participants responded, they received feedback on their performance:  
152 the fixation cross turned green (red) in case of a correct (incorrect) answer. Throughout the  
153 experiment the condition blocks were alternated, the first block being the SD condition for all  
154 participants. Before starting the first block, participants performed one training block per  
155 condition. The purpose of this training was 1) to familiarize participants with the experimental  
156 conditions, 2) to initialize the  $\rho$  and  $\theta$  parameters in the QUEST method for the first  
157 experimental block (initial values were respectively  $\rho = 5.4^\circ$  of visual angle and  $\theta = 6^\circ$  of  
158 rotation). All experiments were performed on a cathode ray monitor, positioned 57 cm from  
159 the subject, with a refresh rate of 160 Hz and a resolution of  $1280 \times 1024$  pixels. The  
160 experiment was coded in MATLAB using the Psychophysics Toolbox (Brainard, 1997). The  
161 stimuli were presented in black on a gray background. Throughout the experiment we  
162 recorded EEG signals.

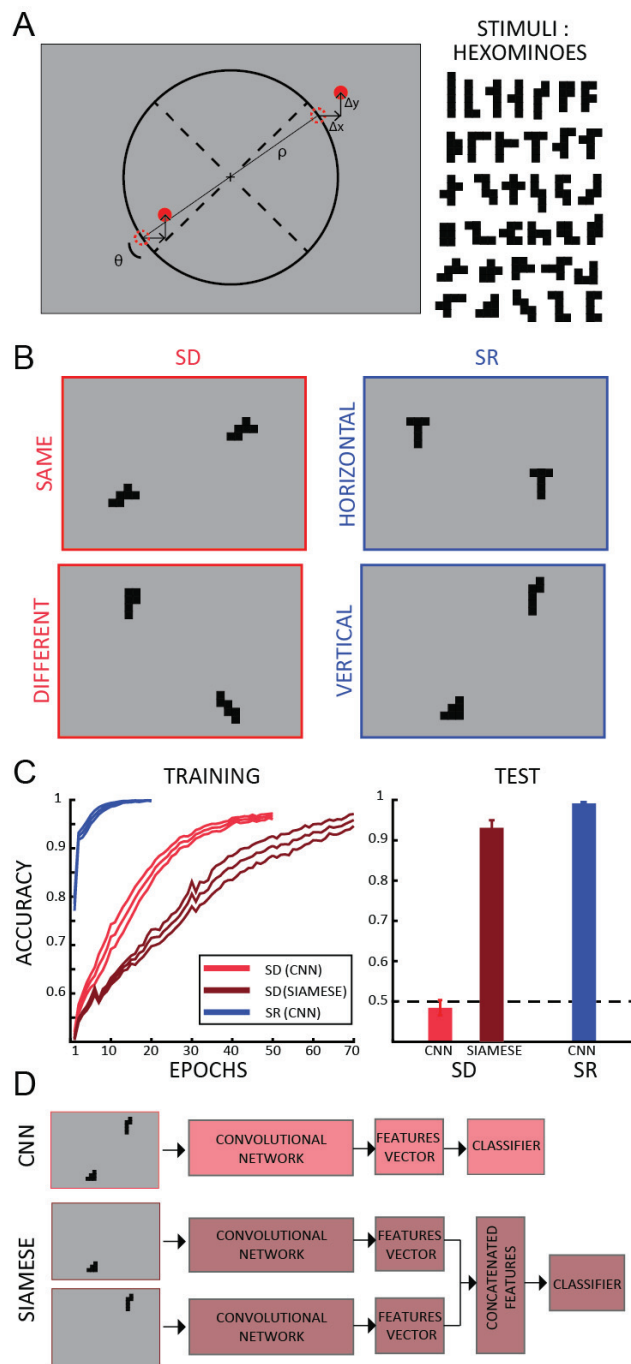
163

### 164 **2.3 EEG recording and pre-processing**

165 We recorded brain activity using a 64-channel active BioSemi electro-  
166 encephalography (EEG) system (1,024 Hz digitizing rate, 3 additional ocular electrodes). The  
167 pre-processing was performed in MATLAB using the EEGLab toolbox (Delorme & Makeig,  
168 2004). First, the data was downsampled to 256 Hz. A notch filter [47Hz - 53Hz] was then  
169 applied to remove power line artifacts. We applied an average-referencing and removed slow  
170 drifts by applying a high-pass filter (>1 Hz). We created the data epochs aligning the data to  
171 the onset of the fixation cross. Finally, we performed an ICA decomposition in order to  
172 remove components related to eye movements and blink artifacts: we visually inspected the  
173 data and removed from 2 to 5 components per subject with a conservative approach (we  
174 removed only components in the frontal regions clearly related to eye movements' activity).

#### 175 **2.4 Computational modeling and code accessibility**

176 We extended a previous computational study (Kim et al., 2018) from which we chose  
177 the parameters of the convolutional feedforward network trained on the SD and SR tasks.  
178 Each task was run 10 times, randomly initializing the networks' parameters and the stimuli  
179 used in the training and test set. The network was fed with 50x80 pixel images. Two  
180 hexominoes (width and height of 2 to 5 pixels) were placed at opposite sides of the screen  
181 (see Fig. 1A and 'Experimental design'). The dictionary of hexominoes was composed of 35  
182 items, which were randomly split between a training (30 items) and a test set (5 items) at  
183 each iteration. Both the training, validation and test sets were composed of 1,000 stimuli (i.e.  
184 different combinations of the hexominoes, with slightly different eccentricity and/or offset  
185 relative to the diagonal). The network consisted of 6 convolutional layers. Each layer  
186 contained 4 channels of size 2x2, with stride of 1. All convolutional layers used a ReLU  
187 activation function with stride of 1 and were followed by pooling layers with 2x2 kernels and a  
188 stride of 1. Eventually, two fully connected layers with 128 units preceded a two-dimensional  
189 classification layer with a sigmoid activation function. As a regularizer we set a dropout rate  
190 of 0.3 in each layer of the network. We used binary cross-entropy as a loss function, the  
191 Adaptive Moment Estimation (Adam) optimizer (Kingma & Ba, 2015) and a learning rate of  
192 10e-4. Each simulation was run over 70 epochs with batch size of 50. All simulations were  
193 run in TensorFlow (GoogleResearch, 2015). The Siamese network had the same exact  
194 convolutional architecture as described above; additionally, the difference between features-  
195 vectors of each separate item (computed on an input image where this item was shown  
196 alone) was fed to the classifier to perform the SD task. All networks count ~7e06 parameters.  
197 All the code and data required to replicate the simulations are available at a github repository  
198 (<https://github.com/artipago/SD-SR>). The code has been run on a Window PC on Python  
199 using the "Tensorflow", "Keras", "Scipy" and "Numpy" libraries.



200

201 **Fig.1: Stimuli and simulation results.** A) The stimuli were the same in the simulations and in the human  
 202 experiments. The items were displayed at opposite sides of the screen (either 45° and 225° or -45° and -225°).  
 203 Both item positions were jittered by a random amount in both the x and y axis ( $\Delta x$  and  $\Delta y$  in the picture) to  
 204 make the task non-trivial for human participants (i.e. preventing participants from performing the SR task  
 205 considering only the position of one item, thus ignoring the spatial relationship between the two items). The  
 206 items used are hexominoes (right panel). Minimum and maximum item height and width are 1.2° – 3.6° and

207 1.2° – 2.7° of visual angle respectively, and 2 to 5 pixels used for the simulations (image size was 50 x 80  
208 pixels).B) Example of stimuli position for the same-different task (SD - left column) and spatial relation task (SR  
209 - right column). For the sake of illustration the ratio between the screen and hexominoes size has been  
210 modified (stimuli here look bigger than in the real experiment). C-D) Accuracy of the CNN network on the  
211 same-different (SD; light red) and spatial relationship (SR; blue) tasks, and of a Siamese network trained on the  
212 SD task (dark red). The Siamese network mimics segmentation in a feedforward network, by separating the  
213 items in two distinct channels of the network (see panel D). The left panel shows the training curves for each  
214 network (accuracy over epochs during training); we stopped the training when the validation accuracy reached  
215 90%. In the right panel we show the training accuracy at the last epoch and the test accuracy. The latter was  
216 evaluated using novel items never used for training, and it reveals that the CNN seems to only learn the  
217 required rule for the SR but not for the SD task, as shown in a previous study. Conversely, the Siamese network  
218 (CNN with segmentation) can solve the SD task, demonstrating that segmentation can allow the CNN to  
219 successfully accomplish this task. In both panels we show average values  $\pm$  SE over 10 repetitions using  
220 different random initializations.

221

## 222 **2.5 Statistical analysis – behavior**

223 We analyzed both accuracy and reaction times (RT) by means of Bayesian ANOVA,  
224 considering the block condition (SR and SD, see above) as independent variables and the  
225 trial condition (whether the stimuli were same or different, or more horizontally or vertically  
226 aligned). The result of such analysis provides a Bayes Factors (BF), which quantifies the  
227 ratio between statistical models given the data. Throughout the paper, all BFs reported  
228 correspond to the probability of the alternative hypothesis over the null hypothesis (indicated  
229 as  $BF_{10}$ ). Practically, a large BF ( $\sim BF > 5$ ) provides evidence in favor of the alternative  
230 hypothesis (the larger the BF the stronger the evidence), whereas low BF ( $\sim BF < 0.5$ )  
231 suggests a lack of effect (Masson, 2011; Smith, 2001). We performed all Bayesian analyses  
232 in JASP (JASP Team, 2018; Love et al., 2015).

## 233 **2.6 Statistical analysis – electrophysiology**

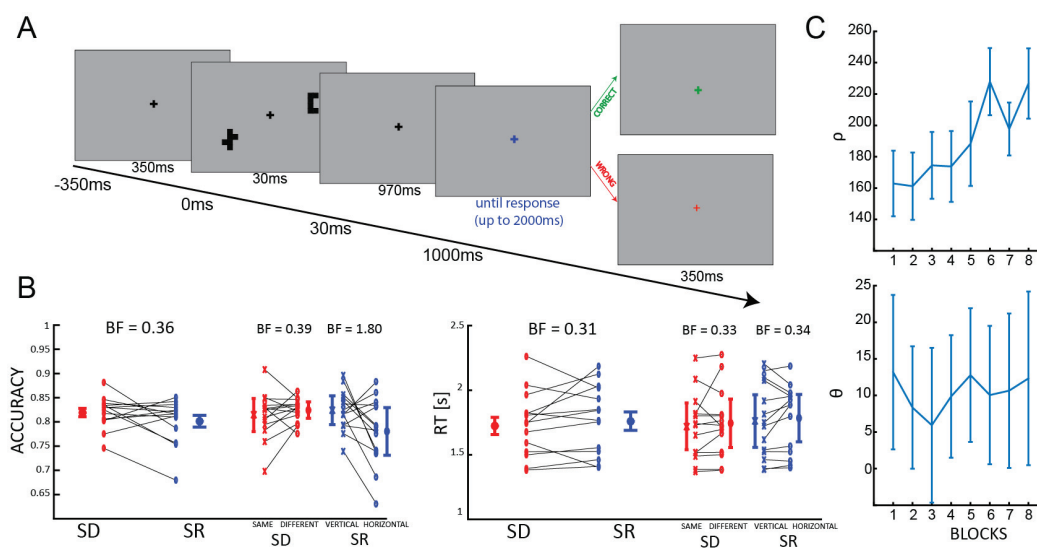
234 Regarding the EEG recording we performed 2 analyses: one in the time domain  
235 measuring Evoked Related Potentials – ERPs, and the other one in the frequency domain  
236 using a time-frequency transform. In the first case, we considered the ERPs recorded from 7  
237 midline electrodes (i.e., Oz, POz, Pz, CPz, Cz, FCz and Fz). After subtracting the baseline  
238 activity recorded during the 350ms before stimuli onset, we averaged the signals from the SD  
239 and SR blocks respectively (i.e., 8 blocks for each condition). Finally, we tested whether the  
240 difference between these signals differed from 0 by means of a point-by-point 2-tailed t test  
241 with a false discovery rate (FDR) correction for multiple comparisons (Hochberg, 1995).  
242 Regarding the time-frequency analysis, we computed the power spectra by means of a  
243 wavelet transform (1–50 Hz in log-space frequency steps with 1-20 cycles). After baseline  
244 correction (i.e., dividing by the averaged activity of the 350ms prior to the onset of the fixation  
245 cross), for each participant, we computed the difference in decibel of the two conditions point  
246 by point, averaging over all electrodes. As in the ERP analysis, we performed a point-by-

247 point 2-tailed t test to identify the time-frequency regions which were significantly different.  
 248 We applied a cluster-based permutation to correct for multiple comparisons (Maris &  
 249 Oostenveld, 2007). First, we identified clusters composed of t values  $t > 3.5$  ( $p < 0.01$ ), and for  
 250 each one we computed the respective global sum. In order to estimate the null distribution  
 251 over the combined t values, we performed the same procedure 500 times after shuffling the  
 252 subject by subject SD-SR assignment. Eventually, we obtained the p values for each non-  
 253 shuffled cluster given the null distribution. All EEG analyses were performed in Matlab; the  
 254 wavelet transform was performed using the EEGLab toolbox (Delorme & Makeig, 2004).

### 255 3. Results

#### 256 3.1 Computational modeling

257 We first extended the results by Kim et al. (2018) for our novel stimulus set: we  
 258 trained two separate Convolutional Neural Networks (CNN) architectures to solve an SD and  
 259 an SR task using a single stimulus set (Methods). The input to these networks was an image  
 260 (50x80 pixels) in which two hexominoes (width and height of 2 to 5 pixels) were displayed at  
 261 opposite sides of the screen (see Fig. 1A). The networks were trained to classify whether the  
 262 two hexominoes were the same or not (SD task) or whether they were aligned more vertically  
 263 or more horizontally with respect to the midline (SR task).  
 264



265

266 **Fig.2: Experimental design and human behavioral results.** A) At the beginning of each trial a black fixation  
 267 cross was displayed for 350ms. After 2 stimuli were shown for 30ms, participants waited an additional 970ms  
 268 before providing the answer. The response was cued by the fixation cross turning blue. After the response, the  
 269 color of the fixation cross provided feedback: green if the response was correct, red otherwise. B) Humans  
 270 performed the SD and SR tasks with comparable levels of performance. In the left and right panels are shown

271 the averages  $\pm$  SE for accuracy and reaction times, respectively. Each pair of connected markers represent an  
272 individual subject. The results for the same-different (in red) and spatial relationship (in blue) conditions are  
273 further broken down for each condition separately (same-different and vertical-horizontal). BF indicates the  
274 Bayes factor against the null hypothesis (difference between the two conditions). C) Changes over blocks of  $p$   
275 (the distance between the stimuli - left panel) and  $\theta$  (the angle between the stimuli and the meridian - right  
276 panel) as adjusted by the QUEST algorithm.

277 We trained and tested the network on different sets of items (a training and test set,  
278 respectively) to assess the networks' ability to generalize beyond training data. We trained  
279 and tested the networks 10 times – randomly initializing networks parameters and training –  
280 test set split each time. We report the mean accuracy and standard deviation over these 10  
281 repetitions in Fig. 1B. Our results are consistent with those from Kim et al. (2018): a CNN  
282 appears to be able to learn the abstract rule (as measured by the network's ability to  
283 generalize beyond the shapes used for training) for SR tasks much more easily than SD  
284 tasks. The effortless ability of humans and other animals (Daniel, Wright, & Katz, 2015;  
285 Wasserman, Castro, & Freeman, 2012) to learn SD tasks suggest the possible involvement  
286 of additional computations that are lacking in CNNs, possibly achieving items identification or  
287 segmentation (e.g. via attention and working memory). In order to verify that segmentation  
288 could be a missing ingredient for the SD task, we implemented a variant of the CNN with  
289 built-in segmentation properties, and tested it on the SD task (it is not necessary to test it on  
290 the SR task, because generalization performance is already at ceiling). The new network  
291 used a Siamese architecture (Bromley et al., 1994) in which each item is processed  
292 separately and eventually combined before being passed to a classifier. Therefore this model  
293 mimics the effect of selective attention and item segregation by feeding to the network each  
294 item separately. The Siamese network could achieve the same training performance on the  
295 SD task as the standard CNN (even though the training took more epochs), however the  
296 network was able to generalize to the test set, while the standard CNN test accuracy was at  
297 chance. This supports the idea that item segmentation or individuation abilities are needed to  
298 achieve the SD task. Next, we test the prediction that SD tasks in humans also require  
299 additional computational mechanisms than SR tasks, by recording EEG signals from a pool  
300 of 28 participants (14 of which were tested on a pilot experiment –fig. 5) performing the same  
301 SD and SR tasks.

### 302 **3.2 Human behavior**

303 A first pilot group of 14 participants performed the SD and SR tasks as described in  
304 Figure 2A, but without any procedure for adjusting task difficulty (i.e. the QUEST method).  
305 The same EEG oscillatory differences between the two tasks as in the main experiment were  
306 observed (fig 5); however, concomitant differences in behavioral task performance left open  
307 the possibility that the oscillatory effects were caused by differences in task difficulty (fig. 5A).

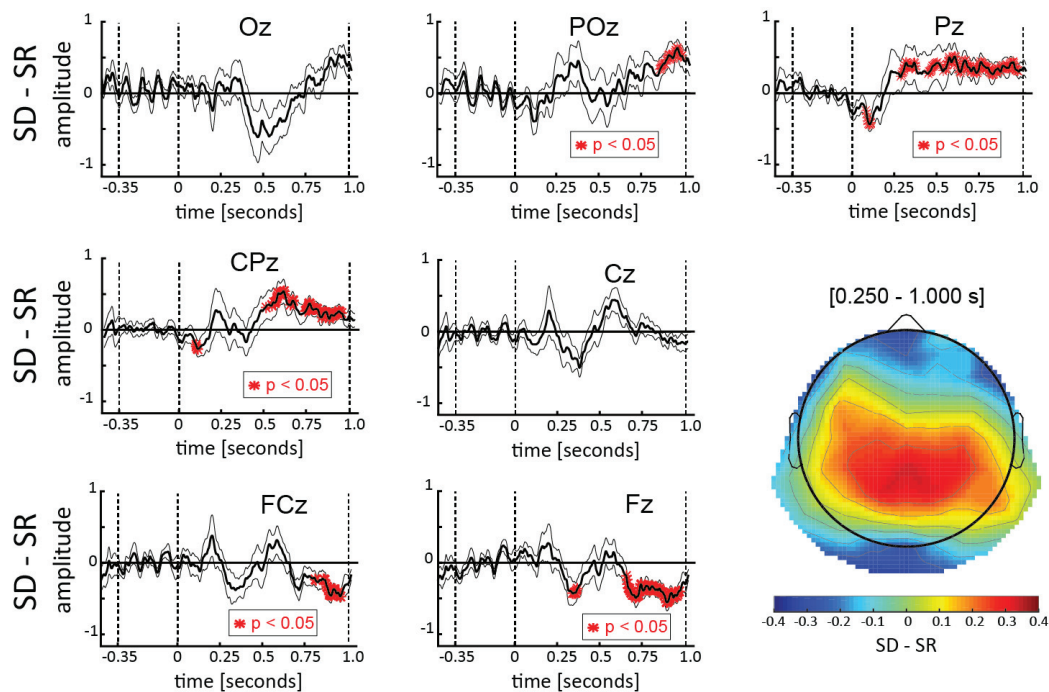
308 Therefore, we replicated the experiment on another group of 14 subjects, this time with an  
309 adaptive procedure to equate behavioral performance between the SD and SR tasks.

310 Participants (N=14) in this main experimental task completed 16 blocks using the  
311 same stimuli as those used to train CNNs (Fig.1): in half of the blocks they were asked to  
312 report whether the two hexominoes were the same or not (SD conditions), in the other half  
313 whether the hexominoes were more vertically or horizontally aligned (SR conditions). The  
314 two conditions were interleaved in a block design. Participants were required to answer after  
315 one second from stimulus onset in order to disentangle motor from visual components in the  
316 EEG recordings (Fig. 2A). The QUEST algorithm was used to assure that participants'  
317 accuracy was matched between the two tasks and remained constant throughout the whole  
318 experiment. This was done by adjusting two experimental parameters trial by trial (i.e., the  
319 hexominoes eccentricity in SD blocks,  $\rho$ , and the angle from the diagonal in SR blocks,  $\theta$ ;  
320 see Fig. 1A and 2C). Maintaining a comparable accuracy between the two tasks reduces the  
321 potential for confounds in the electrophysiological analysis due to differences in performance,  
322 vigilance or motivation. We confirmed the absence of any substantial behavioral difference  
323 between the SD and SR tasks (Fig. 2B) with a Bayesian ANOVA on both accuracy ( $BF_{10} =$   
324  $0.361$ , error  $< 0.001\%$ ) and RT ( $BF_{10} = 0.317$ , error  $< 0.89\%$ ). In addition, we also  
325 investigated each condition separately (Fig. 2B), comparing the difference between 'same'  
326 and 'different' trials (in SD blocks) and 'vertical' and 'horizontal' trials (in SR blocks) in both  
327 RT and accuracy. All comparisons revealed overall no differences between tasks, except for  
328 the accuracy of vertical and horizontal trials in the SR condition, in which the BF proved  
329 inconclusive (accuracy: SD -  $BF_{10} = 0.39$ , error  $< 0.012\%$ ; SR -  $BF_{10} = 1.80$ , error  $< 0.001\%$ ;  
330 RT: SD -  $BF_{10} = 0.333$ , error  $< 0.01\%$ ; SR -  $BF_{10} = 0.34$ , error  $< 0.01\%$ ).

### 331 **3.3 Human electrophysiology: evoked potentials**

332 After having confirmed that performance was equal in the two tasks, we characterized  
333 the evoked potentials (EP) in each task. First, we estimated the difference between SR and  
334 SD conditions considering 7 midline electrodes (Fig.3). The results of a point-by-point t-test  
335 corrected for multiple comparisons revealed a significant difference in central and posterior  
336 electrodes (mostly Pz and CPz) between 250ms after the onset of the stimuli and the  
337 response cue, and the opposite effect in frontal electrodes (FCz and Fz) from 750ms to  
338 1000ms, as confirmed by the topography (Fig.3). Overall these results indicate larger  
339 potentials in visual areas during the SD task than in the SR. Previous studies have shown a  
340 relation between EP amplitude (particularly P300 and late components) with attention  
341 (Itthipuripat, Cha, Byers, & Serences, 2017; Itthipuripat, Cha, Deering, Salazar, & Serences,  
342 2018; Krusemark, Kiehl, & Newman, 2016; Van Voorhis & Hillyard, 1977) and visual working  
343 memory (Fabiani, Karis, & Donchin, 1986; Kok, 2001; McEvoy, 1998). Our results are thus

344 consistent with a larger involvement of executive functions in the SD vs. SR task. In the  
 345 following, we investigated whether this hypothesis is corroborated by corresponding  
 346 oscillatory effects in the time-frequency domain in the main experiment.



347

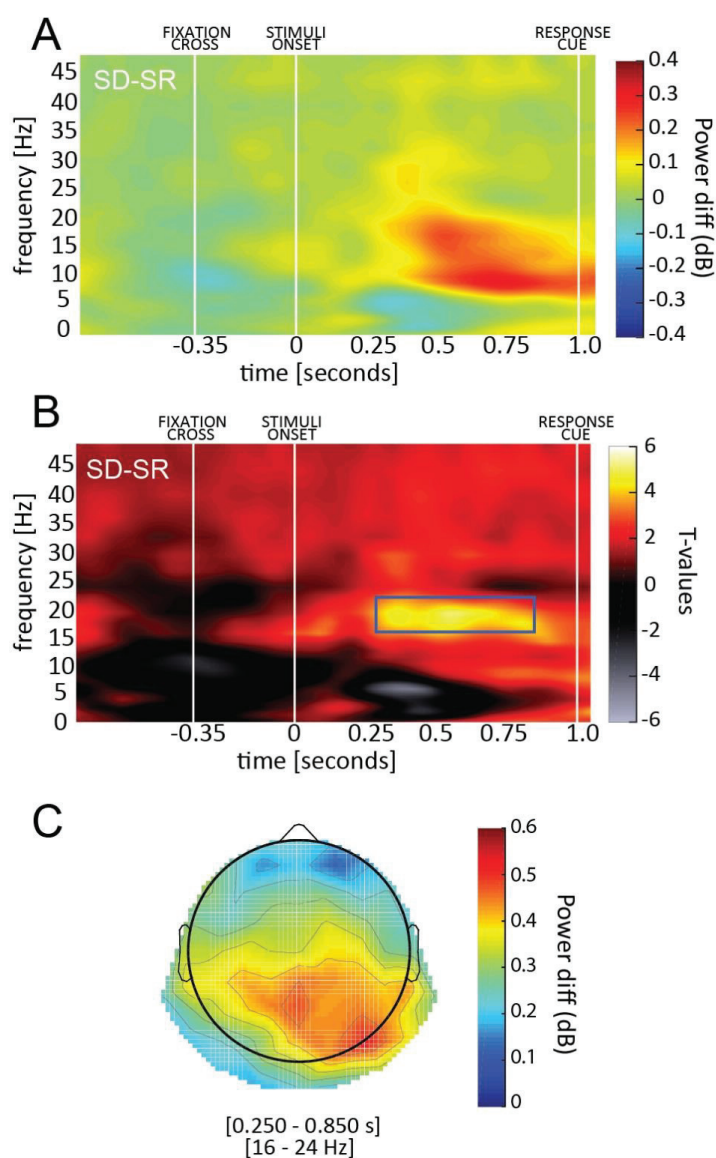
348 **Fig.3: ERPs results.** Each panel represents the difference between ERPs elicited in the SD and SR conditions for the 7 midline electrodes (average  $\pm$  SE). Shown in red are the points for which a significant difference was  
 349 found against zero. The results reveal a significant difference from 250ms after stimuli onset until the response  
 350 cue (at 1000ms) in central parietal regions, and an opposite effect after 750ms in frontal regions. In the  
 351 bottom-right panel the topography, computed over the 250ms – 1000ms interval, confirmed a larger activity in  
 352 the SD than in the SR condition (positive difference, warmer colors) in the central-parietal regions, and an  
 353 opposite effect (negative difference, colder colors) in the frontal regions (which –although not significantly–  
 354 also included occipital regions).  
 355

356

### 357 3.4 Human electrophysiology: time-frequency analysis

358 We performed a time-frequency analysis to try to identify differences between conditions  
 359 observed in specific frequency bands commonly related to executive functions (e.g., visual  
 360 working memory). For this purpose, we computed a baseline-corrected log-scale ratio  
 361 between the two conditions (as shown in Fig. 4A), averaging over all electrodes.  
 362 Remarkably, a point-by-point 2-tailed t-test corrected with cluster-based permutation test  
 363 (Maris & Oostenveld, 2007) revealed a significantly larger activity in the low beta-band (16-  
 364 24Hz) in the SD condition between 250 and 950ms after stimuli onset (Fig. 4B). We further  
 365 quantify the magnitude of the effect by computing the effect size of a one sample t-test

366 against zero averaging per each participant the values within the significant region  
 367 ( $t(13)=2.571$ ,  $p=0.023$ , Cohen's  $d=0.687$ ). The topography of the effect spread mostly over  
 368 parietal and occipital regions (Fig. 4C), mimicking the topography of the EPs analysis. As  
 369 previously, these results confirm the prediction that the SD task may involve additional  
 370 computational mechanisms beyond feedforward computations, possibly indexed by the beta-  
 371 band oscillatory processes identified here. As previously, these results confirmed those from  
 372 the pilot experiment (figure 5D,E), confirming the robustness of the effect also in the  
 373 oscillatory domain. Below, we contextualize and substantiate our results in light of the  
 374 relevant literature.

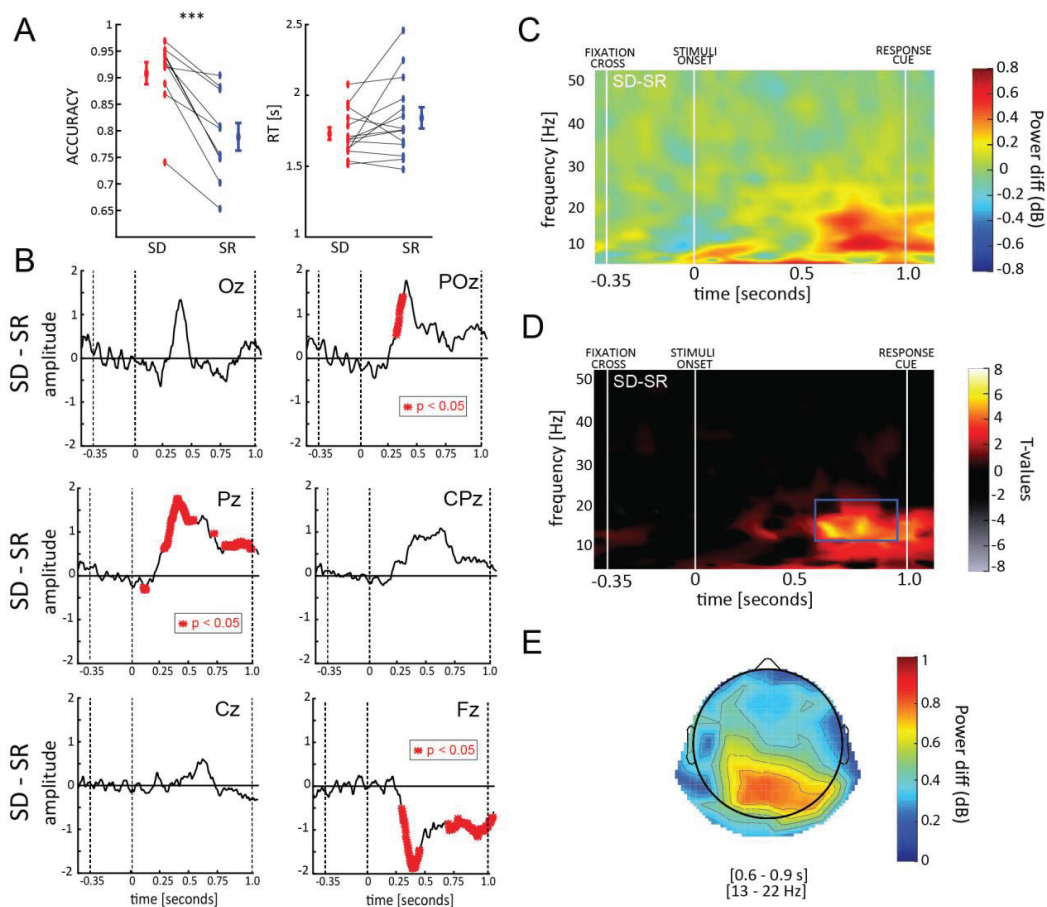


375

376 **Fig.4: Time-frequency results.** A) The difference between SD and SR power spectra is shown in the first panel.  
 377 White lines indicate the onset of the fixation cross, the stimuli and the response cue. B) The second panel  
 378 shows the corresponding t values (when testing the difference against zero). We observed a significant region  
 379 in the low beta band (16-24Hz), between 250ms and 950ms after stimulus onset. C) The topography of the  
 380 significant time-frequency window reveals the involvement of occipital-parietal regions.

381

382



383

384 **Fig.5: Pilot experiment results.** A) Behavioral results of the pilot experiment: left and right panel shows  
 385 accuracy and reaction times for SD (red) and SR (blue) tasks. Differently than in the main task, in the pilot  
 386 experiment participants performed significantly better in the SD than in the SR task (compare the accuracy  
 387 between figure 5A and 2B). B) Difference between SD and SR evoked potentials. Red asterisk indicate time  
 388 window significantly different than zero C) Difference between SD and SR power spectra: white lines indicate  
 389 the stimulus onset and the response cue. D) Testing the SD-SR difference against zero reveals a significant  
 390 region in the low beta band (13-21Hz), before the response cue, in agreement with the results of the main  
 391 experiment –figure 4. We reported a large effect size for this effect (one sample t-test against zero averaging

392 per each participant the values within the significant region,  $t(13)=7.049$ ,  $p<0.001$ , Cohen's  $d=1.820$  E) As in  
393 the main experiment, the SD-SR difference mostly involves occipital-parietal regions.

394

#### 395 **4. Discussion**

396 In this study, we confirmed in a series of two experiments a prediction from the  
397 computational study by Kim et al. (2018) that there exists an important dichotomy between  
398 visual reasoning tasks: While spatial relation (SR) tasks can be solved by modern deep  
399 convolutional neural networks (CNNs), same-different (SD) tasks pose a significant  
400 challenge for these architectures, suggesting the need for additional computations beyond a  
401 feedforward cascade of filtering, rectification and normalization operations. Importantly, the  
402 result of these simulations does not allow us to formulate any prediction about the specific  
403 cortical processes involved in the two tasks. Nonetheless, it demonstrates a fundamental  
404 computational difference, which can be tracked in terms of its human brain neural correlates  
405 while subjects solve SD vs. SR tasks (with difficulty objectively matched by an adaptive  
406 psychometric procedure). Remarkably, in both the pilot and the main experiment we found  
407 higher activity in the former task, in both evoked potentials and oscillatory components. We  
408 interpret these differences as reflecting additional computations required by the SD task. We  
409 can speculate that these additional computations involve working memory and attention  
410 processes, which are lacking in feedforward architectures such as CNNs.

411 Additionally, it is possible to interpret our results in a broader context, by considering  
412 other tasks supposed to involve spatial attention, such as visual search. Previous  
413 experimental work suggested the need for re-entrant processes (Treisman & Gelade, 1980;  
414 Wolfe, Cave, & Franzel, 1989), and how increased activity in specific oscillatory components  
415 (i.e. low [22-34Hz] and high [36-56Hz] gamma bands) are characteristic of these processes  
416 (Buschman & Miller, 2007; S. Phillips & Takeda, 2009). Accordingly, state-of-the-art  
417 computational models performing visual search and related tasks (e.g. instance  
418 segmentation) also employ attentional or recurrent mechanisms (Linsley, Ashok,  
419 Govindarajan, Liu, & Serre, 2020), supporting the hypothesis that convolutional feedforward  
420 networks can benefit from recurrent mechanisms in solving visual reasoning tasks (Kreiman  
421 & Serre, 2020).

422 Computational evidence for the hypothesis that the SD task requires additional  
423 computational mechanisms beyond those needed to solve the SR task is provided by the  
424 results of the Siamese network simulations (Bromley et al., 1994). This feedforward network  
425 processes each stimulus item in a separate (CNN) channel and then passes the processed  
426 items to a single classifier network. Since each item is processed separately (the network is

427 fed two images with only one item represented in each), this 'oracle' architecture performs  
428 the task with item-segmentation processes automatically provided. Our results (as previously  
429 shown on another dataset by Kim and colleagues (Kim et al., 2018) demonstrate that such a  
430 feedforward network, once endowed with object individuation using the Siamese  
431 architecture, can easily learn to solve the SD task. In other words, this model simulates the  
432 beneficial effects of attentional selection, individuation and working memory by segregating  
433 the representations of each item. Our EEG results are compatible with this interpretation,  
434 with higher activity in the SD compared to the SR task, visible in both evoked potentials and  
435 oscillatory frequency bands that have been previously related to attention and working  
436 memory (Benchenane et al., 2011; Lundqvist et al., 2018; Nash & Fernandez, 1996).

437 Previous work has shown that modulations of beta-band oscillations can be related to  
438 selective attention mechanisms (Benchenane et al., 2011; Buschman & Miller, 2007; Lee,  
439 Whittington, & Kopell, 2013; Richter, Coppola, & Bressler, 2018). Different attentional  
440 mechanisms may indeed be involved in the two tasks: the SR task could be solved by first  
441 grouping items and then determining the orientation of the group (Franconeri, Scimeca, Roth,  
442 Helseth, & Kahn, 2012), whereas the SD task requires the individuation of the two items  
443 before comparison. In addition, our results are also consistent with differences in memory  
444 processes between the two tasks (de Fockert, G., Frith, & Lavie, 2001). One common  
445 assumption is that items that are grouped together (as in the SR task) occupy only one  
446 working memory slot (Clevenger & Hummel, 2014; Franconeri, Alvarez, & Cavanagh, 2013),  
447 whereas non-grouped items would each hold one slot, resulting in a larger working memory  
448 load. Previous literature showed that working memory can also be characterized by neuronal  
449 oscillatory signatures. Recent studies, for example, have demonstrated an interplay between  
450 beta and gamma band frequencies during working memory tasks (Lundqvist et al., 2016,  
451 2018). Similarly, alpha and low beta bands, not only increase with working memory load  
452 (Babiloni et al., 2004; Pesonen, Hämäläinen, & Krause, 2007), but also in conjunction with  
453 the inhibition of competing visual memories in selective memory retrieval (Park, Min, & Lee,  
454 2010; Waldhauser, Johansson, & Hanslmayr, 2012). Besides, previous studies have  
455 reported that increased oscillatory activity in the alpha band is a signature of attentional  
456 processes, and it can predict the likelihood of successful trials in many tasks (Händel,  
457 Haarmeier, & Jensen, 2011; Klimesch, 2012; Nelli, Itthipuripat, Srinivasan, & Serences,  
458 2017); however, in our current study we did not investigate differences between correct and  
459 incorrect trials, but between different types of tasks (involving spatial relationship or  
460 sameness judgment), after controlling for task difficulty. This could explain why alpha-band  
461 amplitude differences were less prominent in our study. All considered, several lines of  
462 evidence point towards beta oscillations as crucially involved in both attention and working

463 memory related processes. These processes, therefore, might be part of the additional  
 464 computational mechanisms required for SD tasks compared to SR tasks. Future work could  
 465 more directly compare the attention and memory dependence of each task in human  
 466 subjects.

467 That feedforward neural networks are limited in their ability to solve simple visual  
 468 reasoning tasks is already being recognized by the computer vision and neuroscience  
 469 communities (Kar, Kubilius, Schmidt, Issa, & DiCarlo, 2019; Rajalingham, Issa, Schmidt, Kar,  
 470 & DiCarlo, 2017; Yamins, Hong, Cadieu, & DiCarlo, 2013). Current CNN extensions include  
 471 modules for integrating local and global features (Chen et al., 2018) as well as recurrent  
 472 neural architectures (Yang et al., 2018). Our results suggest that the human visual system  
 473 also deploys additional computations beyond feedforward processes to successfully solve  
 474 visual reasoning tasks. Rhythmic cortical oscillations in the beta-band represent the  
 475 signatures of these additional computations, which may involve selective attention and  
 476 working memory.

#### 477 **Acknowledgments**

478 This work was funded by an ERC Consolidator Grant P-CYCLES number 614244 to RV, a  
 479 joint CRCNS ANR-NSF Grant “OsciDeep” to RV (ANR-19-NEUC-0004) and TS (IIS-  
 480 1912280), and two ANITI (Artificial and Natural Intelligence Toulouse Institute, ANR-19-PI3A-  
 481 0004) Research Chairs to RV and TS. Additional support to TS was provided by ONR grant  
 482 (N00014-19-1-2029).

#### 483 **References**

- 484 Babiloni, C., Babiloni, F., Carducci, F., Cappa, S. F., Cincotti, F., Del Percio, C., ... Rossini, P. M. (2004).  
 485 Human cortical rhythms during visual delayed choice reaction time tasks: A high-resolution EEG study on  
 486 normal aging. *Behavioural Brain Research*, *153*(1), 261–271. <http://doi.org/10.1016/j.bbr.2003.12.012>  
 487 Benchenane, K., Tiesinga, P. H., & Battaglia, F. P. (2011). Oscillations in the prefrontal cortex: A gateway to  
 488 memory and attention. *Current Opinion in Neurobiology*. <http://doi.org/10.1016/j.conb.2011.01.004>  
 489 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, *10*, 433–436.  
 490 <http://doi.org/10.1163/156856897X00357>  
 491 Bromley, J., Guyon, I., Lecun, Y., Sackinger, E., Shah, R., Bell, A., & Holmdel, L. (1994). Signature  
 492 Verification using a “Siamese” Time Delay Neural Network. In *Advances in neural*  
 493 *information processing systems* (pp. 737--744). Retrieved from [http://papers.nips.cc/paper/769-signature-](http://papers.nips.cc/paper/769-signature-verification-using-a-siamese-time-delay-neural-network.pdf)  
 494 [verification-using-a-siamese-time-delay-neural-network.pdf](http://papers.nips.cc/paper/769-signature-verification-using-a-siamese-time-delay-neural-network.pdf)  
 495 Buschman, T. J., & Miller, E. K. (2007). Top-down versus bottom-up control of attention in the prefrontal and  
 496 posterior parietal cortices. *Science*, *315*(5820), 1860–1864. <http://doi.org/10.1126/science.1138071>  
 497 Chen, X., Li, L.-J., Fei-Fei, L., & Gupta, A. (2018). Iterative Visual Reasoning Beyond Convolutions.  
 498 <http://doi.org/10.1109/CVPR.2018.00756>  
 499 Clevenger, P. E., & Hummel, J. E. (2014). Working memory for relations among objects. *Attention, Perception,*  
 500 *and Psychophysics*, *76*(7), 1933–1953. <http://doi.org/10.3758/s13414-013-0601-3>  
 501 Daniel, T. A., Wright, A. A., & Katz, J. S. (2015). Abstract-concept learning of difference in pigeons. *Animal*  
 502 *Cognition*, *18*(4), 831–837. <http://doi.org/10.1007/s10071-015-0849-1>  
 503 de Fockert, J. W., G., R., Frith, C. D., & Lavie, N. (2001). The role of working memory in visual selective  
 504 attention. *S*, *291*, 1803–1806.  
 505 Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-trial EEG dynamics

- 506 including independent component analysis. *Journal of Neuroscience Methods*, 134(1), 9–21.  
 507 <http://doi.org/10.1016/j.jneumeth.2003.10.009>
- 508 Fabiani, M., Karis, D., & Donchin, E. (1986). P300 and Recall in an Incidental Memory Paradigm.  
 509 *Psychophysiology*, 23(3), 298–308. <http://doi.org/10.1111/j.1469-8986.1986.tb00636.x>
- 510 Franconeri, S. L., Alvarez, G. A., & Cavanagh, P. (2013). Flexible cognitive resources: Competitive content  
 511 maps for attention and memory. *Trends in Cognitive Sciences*. <http://doi.org/10.1016/j.tics.2013.01.010>
- 512 Franconeri, S. L., Scimeca, J. M., Roth, J. C., Helseth, S. A., & Kahn, L. E. (2012). Flexible visual processing of  
 513 spatial relationships. *Cognition*, 122(2), 210–227. <http://doi.org/10.1016/j.cognition.2011.11.002>
- 514 GoogleResearch. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. *Google*  
 515 *Research*. <http://doi.org/10.1207/s15326985sep4001>
- 516 Händel, B. F., Haarmeier, T., & Jensen, O. (2011). Alpha oscillations correlate with the successful inhibition of  
 517 unattended stimuli. *Journal of Cognitive Neuroscience*, 23(9), 2494–2502.  
 518 <http://doi.org/10.1162/jocn.2010.21557>
- 519 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the*  
 520 *IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (Vol. 2016–Decem, pp.  
 521 770–778). <http://doi.org/10.1109/CVPR.2016.90>
- 522 Hochberg, B. (1995). Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple  
 523 Testing. *Journal of the Royal Statistical Society*, 57(1), 289–300. <http://doi.org/10.2307/2346101>
- 524 Hollard, V. D., & Delius, J. D. (1982). Rotational invariance in visual pattern recognition by pigeons and  
 525 humans. *Science*, 218(4574), 804–806. <http://doi.org/10.1126/science.7134976>
- 526 Itthipuripat, S., Cha, K., Byers, A., & Serences, J. T. (2017). Two different mechanisms support selective  
 527 attention at different phases of training. *PLoS Biology*, 15(6). <http://doi.org/10.1371/journal.pbio.2001724>
- 528 Itthipuripat, S., Cha, K., Deering, S., Salazar, A. M., & Serences, J. T. (2018). Having more choices changes  
 529 how human observers weight stable sensory evidence. *Journal of Neuroscience*, 38(40), 8635–8649.  
 530 <http://doi.org/10.1523/JNEUROSCI.0440-18.2018>
- 531 JASP Team. (2018). JASP (Version 0.8.6.0). [Computer Software]. Retrieved from <http://jasp-stats.org>
- 532 Kar, K., Kubilius, J., Schmidt, K., Issa, E. B., & DiCarlo, J. J. (2019). Evidence that recurrent circuits are critical  
 533 to the ventral stream’s execution of core object recognition behavior. *Nature Neuroscience*, 22(6), 974–  
 534 983. <http://doi.org/10.1038/s41593-019-0392-5>
- 535 Kim, J., Ricci, M., & Serre, T. (2018). Not-So-CLEVR: Learning same-different relations strains feedforward  
 536 neural networks. *Interface Focus*, 8(4). <http://doi.org/10.1098/rsfs.2018.0011>
- 537 Kingma, D. P., & Ba, J. L. (2015). Adam: A method for stochastic gradient descent. *ICLR: International*  
 538 *Conference on Learning Representations*.
- 539 Klimesch, W. (2012). Alpha-band oscillations, attention, and controlled access to stored information. *Trends in*  
 540 *Cognitive Sciences*. <http://doi.org/10.1016/j.tics.2012.10.007>
- 541 Kok, A. (2001). On the utility of P3 amplitude as a measure of processing capacity. *Psychophysiology*, 38(3),  
 542 557–577. <http://doi.org/10.1017/S0048577201990559>
- 543 Kreiman, G., & Serre, T. (2020). Beyond the feedforward sweep: feedback computations in the visual cortex.  
 544 *Annals of the New York Academy of Sciences*. <http://doi.org/10.1111/nyas.14320>
- 545 Krusemark, E. A., Kiehl, K. A., & Newman, J. P. (2016). Endogenous attention modulates early selective  
 546 attention in psychopathy: An ERP investigation. *Cognitive, Affective and Behavioral Neuroscience*, 16(5),  
 547 779–788. <http://doi.org/10.3758/s13415-016-0430-7>
- 548 Lee, J. H., Whittington, M. A., & Kopell, N. J. (2013). Top-Down Beta Rhythms Support Selective Attention via  
 549 Interlaminar Interaction: A Model. *PLoS Computational Biology*, 9(8).  
 550 <http://doi.org/10.1371/journal.pcbi.1003164>
- 551 Linsley, D., Ashok, A. K., Govindarajan, L. N., Liu, R., & Serre, T. (2020). Stable and expressive recurrent  
 552 vision models. Retrieved from <http://arxiv.org/abs/2005.11362>
- 553 Love, J., Selker, R., Verhagen, J., Marsman, M., Gronau, Q. F., Jamil, T., ... Rouder, J. N. (2015). Software to  
 554 sharpen your stats. *APS Observer*, 28(3), 27–29.
- 555 Lundqvist, M., Herman, P., Warden, M. R., Brincat, S. L., & Miller, E. K. (2018). Gamma and beta bursts  
 556 during working memory readout suggest roles in its volitional control. *Nature Communications*, 9(1).  
 557 <http://doi.org/10.1038/s41467-017-02791-8>
- 558 Lundqvist, M., Rose, J., Herman, P., Brincat, S. L. L., Buschman, T. J. J., & Miller, E. K. K. (2016). Gamma and  
 559 Beta Bursts Underlie Working Memory. *Neuron*, 90(1), 152–164.  
 560 <http://doi.org/10.1016/j.neuron.2016.02.028>
- 561 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data. *Journal of*  
 562 *Neuroscience Methods*, 164(1), 177–190. <http://doi.org/10.1016/j.jneumeth.2007.03.024>
- 563 Masson, M. E. J. (2011). A tutorial on a practical Bayesian alternative to null-hypothesis significance testing,  
 564 679–690. <http://doi.org/10.3758/s13428-010-0049-5>
- 565 McEvoy, L. (1998). Dynamic cortical networks of verbal and spatial working memory: effects of memory load

- 566 and task practice. *Cerebral Cortex*, 8(7), 563–574. <http://doi.org/10.1093/cercor/8.7.563>
- 567 Nash, A. J., & Fernandez, M. (1996). P300 and allocation of attention in dual-tasks. *International Journal of*
- 568 *Psychophysiology*, 23(3), 171–180. [http://doi.org/10.1016/S0167-8760\(96\)00049-9](http://doi.org/10.1016/S0167-8760(96)00049-9)
- 569 Nelli, S., Itthipuripat, S., Srinivasan, R., & Serences, J. T. (2017). Fluctuations in instantaneous frequency
- 570 predict alpha amplitude during visual perception. *Nature Communications*, 8(1).
- 571 <http://doi.org/10.1038/s41467-017-02176-x>
- 572 Park, H. D., Min, B. K., & Lee, K. M. (2010). EEG oscillations reflect visual short-term memory processes for
- 573 the change detection in human faces. *NeuroImage*, 53(2), 629–637.
- 574 <http://doi.org/10.1016/j.neuroimage.2010.06.057>
- 575 Pesonen, M., Hämäläinen, H., & Krause, C. M. (2007). Brain oscillatory 4-30 Hz responses during a visual n-
- 576 back memory task with varying memory load. *Brain Research*, 1138(1), 171–177.
- 577 <http://doi.org/10.1016/j.brainres.2006.12.076>
- 578 Phillips, P. J., Yates, A. N., Hu, Y., Hahn, C. A., Noyes, E., Jackson, K., ... O'Toole, A. J. (2018). Face
- 579 recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms.
- 580 *Proceedings of the National Academy of Sciences of the United States of America*, 115(24), 6171–6176.
- 581 <http://doi.org/10.1073/pnas.1721355115>
- 582 Phillips, S., & Takeda, Y. (2009). Greater frontal-parietal synchrony at low gamma-band frequencies for
- 583 inefficient than efficient visual search in human EEG. *International Journal of Psychophysiology*, 73(3),
- 584 350–354. <http://doi.org/10.1016/j.ijpsycho.2009.05.011>
- 585 Rajalingham, R., Issa, E. B., Schmidt, K., Kar, K., & DiCarlo, J. J. (2017). Feedforward deep neural networks
- 586 diverge from humans and monkeys on core visual object recognition behavior. In *Annual Conference on*
- 587 *Cognitive Computational Neuroscience* (pp. 1–2). Retrieved from
- 588 <https://www2.securecms.com/CCNeuro/docs-0/59288a3f68ed3f3c458a257f.pdf>
- 589 Richter, C. G., Coppola, R., & Bressler, S. L. (2018). Top-down beta oscillatory signaling conveys behavioral
- 590 context in early visual cortex. *Scientific Reports*, 8(1). <http://doi.org/10.1038/s41598-018-25267-1>
- 591 Serre, T. (2019). Deep Learning: The Good, the Bad, and the Ugly. *Annual Review of Vision Science*, 5(1), 399–
- 592 426. <http://doi.org/10.1146/annurev-vision-091718-014951>
- 593 Smith, J. M. B. and A. F. M. (2001). Bayesian Theory. *Measurement Science and Technology*.
- 594 <http://doi.org/10.1088/0957-0233/12/2/702>
- 595 Stabinger, S., Rodríguez-Sánchez, A., & Piater, J. (2016). 25 years of CNNS: Can we compare to human
- 596 abstraction capabilities? In *Lecture Notes in Computer Science (including subseries Lecture Notes in*
- 597 *Artificial Intelligence and Lecture Notes in Bioinformatics)* (Vol. 9887 LNCS, pp. 380–387).
- 598 [http://doi.org/10.1007/978-3-319-44781-0\\_45](http://doi.org/10.1007/978-3-319-44781-0_45)
- 599 Treisman, A., & Gelade, G. (1980). A feature integration theory of attention. *Cognitive Psychology*, 12(1), 97–
- 600 136. Retrieved from <http://www.erc.caltech.edu/Industry/Conferences/2004/AIC/pdf/Laurent-Itti.pdf>
- 601 Van Voorhis, S., & Hillyard, S. A. (1977). Visual evoked potentials and selective attention to points in space.
- 602 *Perception & Psychophysics*, 22(1), 54–62. <http://doi.org/10.3758/BF03206080>
- 603 VanRullen, R., & Thorpe, S. J. (2001). Is it a bird? Is it a plane? Ultra-rapid visual categorisation of natural and
- 604 artificial objects. *Perception*, 30(6), 655–668. <http://doi.org/10.1068/p3029>
- 605 Vogels, R. (1999). Categorization of complex visual images by rhesus monkeys. Part 1: behavioural study.
- 606 *European Journal of Neuroscience*, 11(4), 1223–1238. Retrieved from
- 607 [http://onlinelibrary.wiley.com/doi/10.1046/j.1460-](http://onlinelibrary.wiley.com/doi/10.1046/j.1460-9568.1999.00530.x/full%5Cnhttp://www.ncbi.nlm.nih.gov/pubmed/10103118)
- 608 [9568.1999.00530.x/full%5Cnhttp://www.ncbi.nlm.nih.gov/pubmed/10103118](http://www.ncbi.nlm.nih.gov/pubmed/10103118)
- 609 Waldhauser, G. T., Johansson, M., & Hanslmayr, S. (2012). Alpha/Beta Oscillations Indicate Inhibition of
- 610 Interfering Visual Memories. *Journal of Neuroscience*, 32(6), 1953–1961.
- 611 <http://doi.org/10.1523/jneurosci.4201-11.2012>
- 612 Wasserman, E. A., Castro, L., & Freeman, J. H. (2012). Same-different categorization in rats. *Learning and*
- 613 *Memory*, 19(4), 142–145. <http://doi.org/10.1101/lm.025437.111>
- 614 Watson, A. B., & Pelli, D. G. (1983). Quest: A Bayesian adaptive psychometric method. *Perception &*
- 615 *Psychophysics*, 33(2), 113–120. <http://doi.org/10.3758/BF03202828>
- 616 Wolfe, J. M., Cave, K. R., & Franzel, S. L. (1989). Guided Search: An Alternative to the Feature Integration
- 617 Model for Visual Search. *Journal of Experimental Psychology: Human Perception and Performance*,
- 618 15(3), 419–433. <http://doi.org/10.1037/0096-1523.15.3.419>
- 619 Yamins, D., Hong, H., Cadieu, C., & DiCarlo, J. J. (2013). Hierarchical modular optimization of convolutional
- 620 networks achieves representations similar to Macaque IT and human ventral stream. In *Advances in Neural*
- 621 *Information Processing Systems*.
- 622 Yang, G. R., Ganichev, I., Wang, X. J., Shlens, J., & Sussillo, D. (2018). A dataset and architecture for visual
- 623 reasoning with a working memory. In *Lecture Notes in Computer Science (including subseries Lecture*
- 624 *Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* (Vol. 11214 LNCS, pp. 729–745).
- 625 [http://doi.org/10.1007/978-3-030-01249-6\\_44](http://doi.org/10.1007/978-3-030-01249-6_44)