

Cognition and Behavior

Differential Involvement of EEG Oscillatory Components in Sameness versus Spatial-Relation Visual Reasoning Tasks

[©]Andrea Alamia,¹ Canhuang Luo,¹ Matthew Ricci,² Junkyung Kim,² Thomas Serre,^{2,3} and [©]Rufin VanRullen^{1,3}

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

¹CerCo, Centre National de la Recherche Scientifique Université de Toulouse, Toulouse 31055, France, ²Department of Cognitive, Linguistic and Psychological Sciences, Carney Institute for Brain Science, Brown University, Providence, RI 02912, and ³Artificial and Natural Intelligence Toulouse Institute (ANITI), Université de Toulouse, Toulouse 31055, France

Abstract

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

Key words: deep neural networks; EEG oscillations; ERPs; spatial relationship; visual reasoning

Significance Statement

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

Received June 15, 2020; accepted October 21, 2020; First published November 25, 2020.

The authors declare no competing financial interests.

Author contributions: A.A., M.R., J.K., T.S., and R.V. designed research; A.A. and C.L. performed research; A.A. and R.V. analyzed data; A.A., C.L., M.R., T.S., and R.V. wrote the paper.



Introduction

The field of artificial vision witnessed an impressive boost in the last few years, driven by the striking results of deep convolutional neural networks (CNNs). Such hierarchical neural networks process information sequentially, through a feedforward cascade of filtering, rectification, and normalization operations. The accuracy of these architectures is now approaching, sometimes exceeding, that of human observers on key visual recognition tasks including object (He et al., 2016) and face recognition (Phillips et al., 2018). These advances suggest that purely feedforward mechanisms suffice to accomplish remarkable results in object categorization, in line with previous experimental studies on humans (VanRullen and Thorpe, 2001) and animals (Hollard and Delius, 1982; Vogels, 1999). However, despite the remarkable accuracy reached in these recognition tasks, the limitations of CNNs are becoming increasingly evident (for recent review, see Serre, 2019). Beyond image categorization tasks, CNNs appear to struggle to learn to solve relatively simple visual reasoning tasks otherwise trivial for the human brain (Stabinger et al., 2016; Kim et al., 2018). A recent study (Kim et al., 2018) thoroughly investigated the ability of CNN architectures to learn to solve various visual reasoning tasks, and found an apparent dichotomy between two sorts of problems: on the one hand, tasks that require judging the spatial relations between items [spatial relationship (SR)]; on the other, those that require comparing items [same-different (SD)]. Importantly, Kim and colleagues demonstrated that CNNs can more easily learn the first class of problems compared with the second one.

This prompts the question of how biological visual systems handle such tasks so efficiently. Kim et al. (2018) suggest that SR and SD tasks tap into distinct computational mechanisms, thus leading to the prediction that different cortical processes are also involved when humans perform the two tasks: SR tasks can be successfully solved by feedforward processes, whereas SD tasks seem to require additional computations, such as working memory and attention. Here, we tested this hypothesis in two steps: first, we confirmed and extended Kim's results by comparing the performance of CNNs on an experiment in which we directly contrasted SD and SR tasks on the same stimulus set. Second, we recorded electrophysiological responses [electroencephalography (EEG)] in healthy human participants for the

This work was supported by the ERC Consolidator Grant P-CYCLES 614244 (to R.V.), a Joint Collaborative Research in Computational NeuroScience (CRCNS) Agence Nationale Recherche-National Science Fondation (ANR-NSF) Grant "OsciDeep" to R.V. (ANR-19-NEUC-0004) and T.S. (IIS-1912280), and two Artificial and Natural Intelligence Toulouse Institute (ANITI) ANR-19-PI3A-0004) Research Chairs (to R.V. and T.S.). Additional support to T.S. was provided by the Office of Naval Research (ONR) Grant N00014-19-1-2029.

Correspondence should be addressed to Rufin VanRullen at rufin.vanrullen@cnrs.fr.

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

Copyright © 2021 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.

same experiment, after having matched the difficulty level via an adaptive psychometric procedure. We hypothesized that the additional computations required by the SD task, as compared with SR tasks, would elicit differences in evoked potentials (EPs; e.g., P300 modulations, which have been related to attentional mechanisms; Nash and Fernandez, 1996) and brain rhythms related to working memory (such as β band oscillations; Benchenane et al., 2011; Lundqvist et al., 2018). We found indeed that, in addition to a variation in EPs, the SD task elicited higher activity in specific β band oscillatory components in the occipital-parietal areas, which are typically associated with attention-related and memory-related processes. We emphasize that the goal of the present study was not to identify the precise neural computations involved in the two tasks (which would naturally require a broader experimental set-up than a single EEG study), but rather to validate the hypothesis that SD involves additional computations relative to SR (even when the two tasks are equally difficult). We hope that this demonstration can be a first step toward characterizing the processes taking place in visual cortex during visual reasoning tasks, and designing more reliable and more human-like computational models.

Materials and Methods

Participants and pilot experiment

Twenty-eight participants (aged 21-34 years old with a mean age of 26.6 ± 3.7, 11 women, five left-handed), volunteered to join the experiment. All subjects reported normal or corrected to normal vision and had no history of epileptic seizures or neurologic disorders. Participants were pooled in two groups of 14 each: one group performed a pilot experiment, while the second one was tested on a final version of the task. The only difference between the pilot and the main study was the QUEST adaptive procedure used to match the difficulty level between conditions, which was not implemented in the pilot experiment. However, in both studies we found the very same result (see below, Figs. 3, 4 and 5). In the main experiment, we kept the same number of participants to replicate the effect, after having removed the behavioral difference in task difficulty via the QUEST algorithm. This study complies with the guidelines of the research center where it was conducted, and the protocol was approved by an external committee (ethics approval number N° 2016-A01937-44). All participants gave written informed consent before starting the experiment, in accordance with the Declaration of Helsinki, and received monetary compensation for their participation.

Experimental design

The experiment was composed of 16 experimental blocks of 70 trials each, with a total duration of \sim 1 h. Each trial lasted \sim 2 s (Fig. 1A): 350 ms after the onset of a black fixation cross (0.6° width), two shapes were displayed for 30 ms on opposite sides of the screen, distant $2^*\rho$ from each other with an angle of \pm (45° + θ) with respect to the horizontal midline (ρ being the distance from the center of the screen, and θ the angular difference with



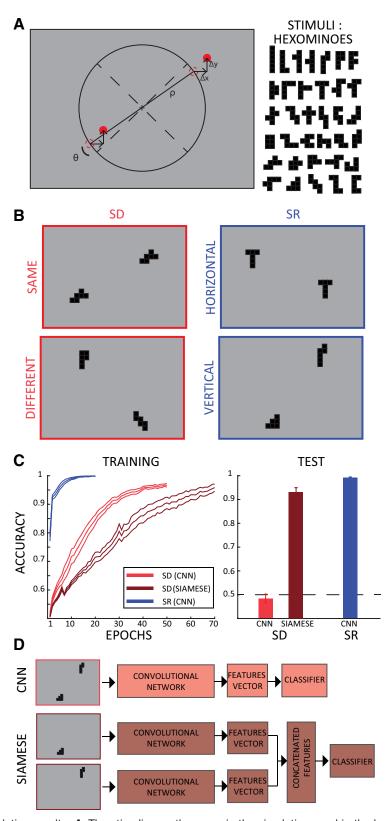


Figure 1. Stimuli and simulation results. A, The stimuli were the same in the simulations and in the human experiments. The items were displayed at opposite sides of the screen (either 45° and 225° or -45° and -225°). Both item positions were jittered by a random amount in both the x- and y-axes (Δx and Δy in the picture) to make the task non-trivial for human participants (i.e., preventing participants from performing the SR task considering only the position of one item, thus ignoring the SR between the two items). The items used are hexominoes (right panel). Minimum and maximum item height and width are 1.2–3.6° and 1.2–2.7° of visual angle, respectively, and 2–5 pixels used for the simulations (image size was 50×80 pixels). $\textbf{\textit{B}}$, Example of stimuli position for the



continued

SD task (left column) and spatial relation task (SR, right column). For the sake of illustration, the ratio between the screen and hexominoes size has been modified (stimuli here look bigger than in the real experiment). C, D, Accuracy of the CNN network on the SD (light red) and SR (blue) tasks, and of a Siamese network trained on the SD task (dark red). The Siamese network mimics segmentation in a feedforward network, by separating the items in two distinct channels of the network (see D). The left panel shows the training curves for each network (accuracy over epochs during training); we stopped the training when the validation accuracy reached 90%. In the right panel, we show the training accuracy at the last epoch and the test accuracy. The latter was evaluated using novel items never used for training, and it reveals that the CNN seems to only learn the required rule for the SR but not for the SD task, as shown in a previous study. Conversely, the Siamese network (CNN with segmentation) can solve the SD task, demonstrating that segmentation can allow the CNN to successfully accomplish this task. In both panels we show average values \pm SE over 10 repetitions using different random initializations.

the diagonal; see Fig. 1B). Each shape was selected from a subset of 36 hexominoes, a geometric figure composed of six contiguous squares (Fig. 1B) One second after the onset of the hexominoes, the fixation cross turned blue, cuing participants to respond. In half of the blocks, participants had to report whether the two shapes were the same or different (SD condition); in the remaining blocks participants had to judge whether the two stimuli were aligned more horizontally or vertically (SR condition). Shapes were displayed at opposite sides of the screen along two main possible orientation axes sampled at random for every trial (either 45° and 225° or -45° and -225°). Both stimuli positions were jittered by a random offset Δx and Δy in both the x and y-axis and a rotation θ from the main axis. The same offsets were applied to both shapes, so they did not affect the angle between stimuli. The aim of such offsets was to prevent participants in the SR condition from determining the configuration of the two stimuli (orientation task) by merely judging the position of a single stimulus: without the random offsets, considering for example the top-right corner position, if the item were below/above the (imaginary) screen diagonal line, the overall orientation would be horizontal/vertical, without the need to consider the position of the corresponding bottom-left item. The offset then compelled participants to consider the relative position of both hexominoes at once. Importantly, in the main experiment (compared with the pilot experiment) the difficulty of the two tasks was controlled by an adaptive psychometric procedure (QUEST method; Watson and Pelli, 1983), which varied the eccentricity of the two stimuli ρ (in the SD blocks) or θ (in the SR blocks) to maintain an overall accuracy level of 80% throughout the whole experiment. In fact, larger (smaller) values of ρ made the stimuli more (less) eccentric and the task more (less) difficult; similarly, smaller (larger) values of θ set the stimuli closer to (farther from) the 45° diagonal line, making the task more (less) difficult. We modified one parameter per condition (i.e., per block), while the other was kept constant (using the same value as in the preceding block). After participants responded, they received feedback on their performance: the fixation cross turned green (red) in case of a correct (incorrect) answer. Throughout the experiment the condition blocks were alternated, the first block being the SD condition for all participants. Before starting the first block, participants performed one training block per condition. The purpose of this training was (1) to familiarize participants with the experimental conditions, (2) to initialize

the ρ and θ parameters in the QUEST method for the first experimental block (initial values were respectively $\rho=5.4^\circ$ of visual angle and $\theta=6^\circ$ of rotation). All experiments were performed on a cathode ray monitor, positioned 57 cm from the subject, with a refresh rate of 160 Hz and a resolution of 1280×1024 pixels. The experiment was coded in MATLAB using the Psychophysics Toolbox (Brainard, 1997). The stimuli were presented in black on a gray background. Throughout the experiment we recorded EEG signals.

EEG recording and preprocessing

We recorded brain activity using a 64-channel active BioSemi EEG system (1024-Hz digitizing rate, three additional ocular electrodes). The preprocessing was performed in MATLAB using the EEGlab toolbox (Delorme and Makeig, 2004). First, the data were downsampled to 256 Hz. A notch filter (47-53 Hz) was then applied to remove power line artifacts. We applied an average-referencing and removed slow drifts by applying a high-pass filter (>1 Hz). We created the data epochs aligning the data to the onset of the fixation cross. Finally, we performed an ICA decomposition to remove components related to eye movements and blink artifacts: we visually inspected the data and removed from two to five components per subject with a conservative approach (we removed only components in the frontal regions clearly related to eye movements' activity).

Computational modeling and code accessibility

We extended a previous computational study (Kim et al., 2018) from which we chose the parameters of the convolutional feedforward network trained on the SD and SR tasks. Each task was run 10 times, randomly initializing the networks' parameters and the stimuli used in the training and test set. The network was fed with 50×80 pixel images. Two hexominoes (width and height of two to five pixels) were placed at opposite sides of the screen (Fig. 1A; see above, Experimental design). The dictionary of hexominoes was composed of 35 items, which were randomly split between a training (30 items) and a test set (five items) at each iteration. Both the training, validation and test sets were composed of 1000 stimuli (i.e., different combinations of the hexominoes, with slightly different eccentricity and/or offset relative to the diagonal). The network consisted of six convolutional layers. Each layer contained four channels of size 2×2 , with stride of 1. All convolutional layers used a ReLu activation function with



stride of 1 and were followed by pooling layers with 2×2 kernels and a stride of 1. Eventually, two fully connected layers with 128 units preceded a two-dimensional classification layer with a sigmoid activation function. As a regularizer we set a dropout rate of 0.3 in each layer of the network. We used binary cross-entropy as a loss function, the adaptive moment estimation (Adam) optimizer (Kingma and Ba, 2015) and a learning rate of 10e-4. Each simulation was run over 70 epochs with batch size of 50. All simulations were run in TensorFlow (GoogleResearch, 2016). The Siamese network had the same exact convolutional architecture as described above; additionally, the difference between features-vectors of each separate item (computed on an input image where this item was shown alone) was fed to the classifier to perform the SD task. All networks count ~7e06 parameters. All the code and data required to replicate the simulations are available at a GitHub repository (https://github. com/artipago/SD-SR). The code has been run on a Window PC on Python using the "Tensorflow," "Keras," "Scipy," and "Numpy" libraries.

Statistical analysis, behavior

We analyzed both accuracy and reaction times (RTs) by means of Bayesian ANOVA, considering the block condition (SR and SD, see above) as independent variables and the trial condition (whether the stimuli were same or different, or more horizontally or vertically aligned). The result of such analysis provides a Bayes factors (BF), which quantifies the ratio between statistical models given the data. Throughout the paper, all BFs reported correspond to the probability of the alternative hypothesis over the null hypothesis (indicated as BF₁₀). Practically, a large BF (\sim BF > 5) provides evidence in favor of the alternative hypothesis (the larger the BF the stronger the evidence), whereas low BF (\sim BF < 0.5) suggests a lack of effect (Bernardo and Smith, 2009; Masson, 2011). We performed all Bayesian analyses in JASP (Love et al., 2015; JASP Team, 2018).

Statistical analysis, electrophysiology

Regarding the EEG recording we performed two analyses: one in the time domain measuring evoked related potentials, ERPs, and the other one in the frequency domain using a time-frequency transform. In the first case, we considered the ERPs recorded from seven midline electrodes (i.e., Oz, POz, Pz, CPz, Cz, FCz, and Fz). After subtracting the baseline activity recorded during the 350 ms before stimuli onset, we averaged the signals from the SD and SR blocks respectively (i.e., eight blocks for each condition). Finally, we tested whether the difference between these signals differed from 0 by means of a pointby-point two-tailed t test with a false discovery rate (FDR) correction for multiple comparisons (Hochberg, 1995). Regarding the time-frequency analysis, we computed the power spectra by means of a wavelet transform (1–50 Hz in log-space frequency steps with 1-20 cycles). After baseline correction (i.e., dividing by the averaged activity of the 350 ms before the onset of the fixation cross), for each participant, we computed the difference in decibel of the two conditions point by point, averaging over all electrodes. As in the ERP analysis, we performed a point-by-point two-tailed t test to identify the time-frequency regions which were significantly different. We applied a cluster-based permutation to correct for multiple comparisons (Maris and Oostenveld, 2007). First, we identified clusters composed of t values t>3.5 (p<0.01), and for each one we computed the respective global sum. In order to estimate the null distribution over the combined t values, we performed the same procedure 500 times after shuffling the subject by subject SD-SR assignment. Eventually, we obtained the p values for each non-shuffled cluster given the null distribution. All EEG analyses were performed in MATLAB; the wavelet transform was performed using the EEGlab toolbox (Delorme and Makeig, 2004).

Results

Computational modeling

We first extended the results by Kim et al. (2018) for our novel stimulus set: we trained two separate CNNs architectures to solve an SD and an SR task using a single stimulus set (Materials and Methods). The input to these networks was an image (50×80 pixels) in which two hexominoes (width and height of two to five pixels) were displayed at opposite sides of the screen (Fig. 1A). The networks were trained to classify whether the two hexominoes were the same or not (SD task) or whether they were aligned more vertically or more horizontally with respect to the midline (SR task).

We trained and tested the network on different sets of items (a training and test set, respectively) to assess the networks' ability to generalize beyond training data. We trained and tested the networks 10 times, randomly initializing networks, parameters as well as the training and test set split each time. We report the mean accuracy and standard deviation over these 10 repetitions in Figure 1B. Our results are consistent with those from Kim et al. (2018): a CNN appears to be able to learn the abstract rule (as measured by the network's ability to generalize beyond the shapes used for training) for SR tasks much more easily than SD tasks. The effortless ability of humans and other animals (Wasserman et al., 2012; Daniel et al., 2015) to learn SD tasks suggest the possible involvement of additional computations that are lacking in CNNs, possibly achieving items identification or segmentation (e.g., via attention and working memory). In order to verify that segmentation could be a missing ingredient for the SD task, we implemented a variant of the CNN with built-in segmentation properties, and tested it on the SD task (it is not necessary to test it on the SR task, because generalization performance is already at ceiling). The new network used a Siamese architecture (Bromley et al., 1993) in which each item is processed separately and eventually combined before being passed to a classifier. Therefore, this model mimics the effect of selective attention and item segregation by feeding to the network each item separately. The Siamese network could achieve the same training performance on the SD task as the standard CNN (although the training took more epochs); however, the



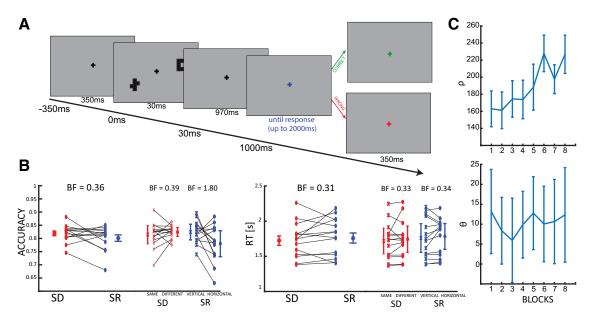


Figure 2. Experimental design and human behavioral results. \bf{A} , At the beginning of each trial, a black fixation cross was displayed for 350 ms. After two stimuli were shown for 30 ms, participants waited an additional 970 ms before providing the answer. The response was cued by the fixation cross turning blue. After the response, the color of the fixation cross provided feedback: green if the response was correct, red otherwise. \bf{B} , Humans performed the SD and SR tasks with comparable levels of performance. In the left and right panels are shown the averages \pm SE for accuracy and RTs, respectively. Each pair of connected markers represent an individual subject. The results for the SD (in red) and SR (in blue) conditions are further broken down for each condition separately (SD and vertical-horizontal). BF indicates the BF against the null hypothesis (difference between the two conditions). \bf{C} , Changes over blocks of ρ (the distance between the stimuli, left panel) and θ (the angle between the stimuli and the meridian, right panel) as adjusted by the QUEST algorithm.

network was able to generalize to the test set, while the standard CNN test accuracy was at chance. This supports the idea that item segmentation or individuation abilities are needed to achieve the SD task Next, we test the prediction that SD tasks in humans also require additional computational mechanisms than SR tasks by recording EEG signals from a pool of 28 participants (14 of which were tested on a pilot experiment; Fig. 5) performing the same SD and SR tasks.

Human behavior

A first pilot group of 14 participants performed the SD and SR tasks as described in Figure 2A, but without any procedure for adjusting task difficulty (i.e., the QUEST method). The same EEG oscillatory differences between the two tasks as in the main experiment were observed (Fig. 5); however, concomitant differences in behavioral task performance left open the possibility that the oscillatory effects were caused by differences in task difficulty (Fig. 5A). Therefore, we replicated the experiment on another group of 14 subjects, this time with an adaptive procedure to equate behavioral performance between the SD and SR tasks.

Participants (*N* = 14) in this main experimental task completed 16 blocks using the same stimuli as those used to train CNNs (Fig. 1): in half of the blocks they were asked to report whether the two hexominoes were the same or not (SD conditions), in the other half whether the hexominoes were more vertically or horizontally aligned (SR conditions). The two conditions were interleaved in a block design. Participants were

required to answer after one second from stimulus onset to disentangle motor from visual components in the EEG recordings (Fig. 2A). The QUEST algorithm was used to assure that participants' accuracy was matched between the two tasks and remained constant throughout the whole experiment. This was done by adjusting two experimental parameters trial by trial (i.e., the hexominoes eccentricity in SD blocks, ρ , and the angle from the diagonal in SR blocks, θ ; see Figs. 1A, 2C). Maintaining a comparable accuracy between the two tasks reduces the potential for confounds in the electrophysiological analysis because of differences in performance, vigilance, or motivation. We confirmed the absence of any substantial behavioral difference between the SD and SR tasks (Fig. 2B) with a Bayesian ANOVA on both accuracy (BF₁₀ = 0.361, error < 0.001%) and RT (BF₁₀ = 0.317, error < 0.89%). In addition, we also investigated each condition separately (Fig. 2B), comparing the difference between "same" and "different" trials (in SD blocks) and "vertical" and "horizontal" trials (in SR blocks) in both RT and accuracy. All comparisons revealed overall no differences between tasks, except for the accuracy of vertical and horizontal trials in the SR condition, in which the BF proved inconclusive (accuracy: SD, BF₁₀ = 0.39, error < 0.012%; SR, BF₁₀ = 1.80, error < 0.001%; RT: SD, BF₁₀ = 0.333, error < 0.01%; SR, $BF_{10} = 0.34$, error < 0.01%).

Human electrophysiology: EPs

After having confirmed that performance was equal in the two tasks, we characterized the EPs in each task.



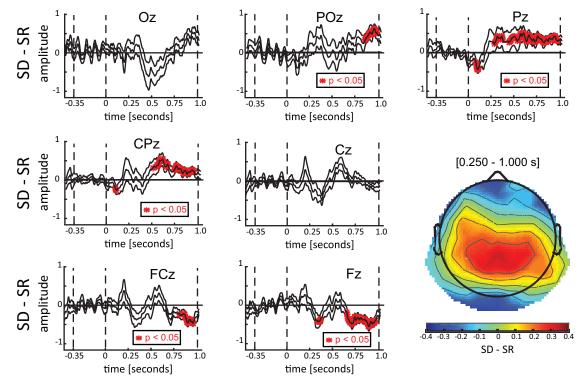


Figure 3. ERPs results. Each panel represents the difference between ERPs elicited in the SD and SR conditions for the seven midline electrodes (average \pm SE). Shown in red are the points for which a significant difference was found against zero. The results reveal a significant difference from 250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

First, we estimated the difference between SR and SD conditions considering 7 midline electrodes (Fig. 3). The results of a point-by-point t test corrected for multiple comparisons revealed a significant difference in central and posterior electrodes (mostly Pz and CPz) between 250 ms after the onset of the stimuli and the response cue, and the opposite effect in frontal electrodes (FCz and Fz) from 750 to 1000 ms, as confirmed by the topography (Fig. 3). Overall, these results indicate larger potentials in visual areas during the SD task than in the SR. Previous studies have shown a relation between EP amplitude (particularly P300 and late components) with attention (Van Voorhis and Hillyard, 1977; Krusemark et al., 2016; Itthipuripat et al., 2017, 2018) and visual working memory (Fabiani et al., 1986; McEvoy et al., 1998; Kok, 2001). Our results are thus consistent with a larger involvement of executive functions in the SD versus SR task. In the following, we investigated whether this hypothesis is corroborated by corresponding oscillatory effects in the time-frequency domain in the main experiment.

Human electrophysiology: time-frequency analysis

We performed a time-frequency analysis to try to identify differences between conditions observed in specific frequency bands commonly related to executive functions (e.g., visual working memory). For this purpose, we computed a baseline-corrected log-scale ratio between the two conditions (as shown in Fig. 4A), averaging over all electrodes. Remarkably, a point-by-point two-tailed t test corrected with cluster-based permutation test (Maris and Oostenveld, 2007) revealed a significantly larger activity in the low β band (16– 24 Hz) in the SD condition between 250 and 950 ms after stimuli onset (Fig. 4B). We further quantify the magnitude of the effect by computing the effect size of a one sample t test against zero averaging per each participant the values within the significant region ($t_{(13)} = 2.571$, p = 0.023, Cohen's d = 0.687). The topography of the effect spread mostly over parietal and occipital regions (Fig. 4C), mimicking the topography of the EPs analysis. As previously, these results confirm the prediction that the SD task may involve additional computational mechanisms beyond feedforward computations, possibly indexed by the β band oscillatory processes identified here. As previously, these results confirmed those from the pilot experiment (Fig. 5D,E), confirming the robustness of the effect also in the oscillatory domain. Below, we contextualize and substantiate our results in light of the relevant literature.

Discussion

In this study, we confirmed in a series of two experiments a prediction from the computational study by Kim et al. (2018) that there exists an important dichotomy



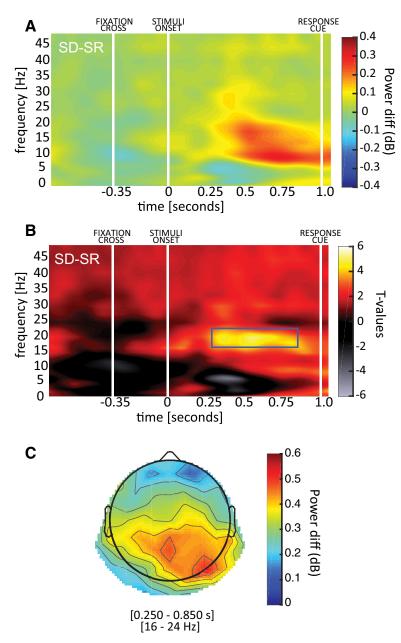


Figure 4. Time-frequency results. **A**, The difference between SD and SR power spectra is shown in the first panel. White lines indicate the onset of the fixation cross, the stimuli and the response cue. **B**, The second panel shows the corresponding t values (when testing the difference against zero).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

between visual reasoning tasks: While SR tasks can be solved by modern deep CNNs, SD tasks pose a significant challenge for these architectures, suggesting the need for additional computations beyond a feedforward cascade of filtering, rectification and normalization operations. Importantly, the result of these simulations does not allow us to formulate any prediction about the specific cortical processes involved in the two tasks. Nonetheless, it demonstrates a fundamental computational difference, which can be tracked in terms of its human brain neural correlates while subjects solve SD versus SR tasks (with difficulty objectively matched by an adaptive psychometric

procedure). Remarkably, in both the pilot and the main experiment we found higher activity in the former task, in both EPs and oscillatory components. We interpret these differences as reflecting additional computations required by the SD task. We can speculate that these additional computations involve working memory and attention processes, which are lacking in feedforward architectures such as CNNs.

Additionally, it is possible to interpret our results in a broader context, by considering other tasks supposed to involve spatial attention, such as visual search. Previous experimental work suggested the need for re-entrant



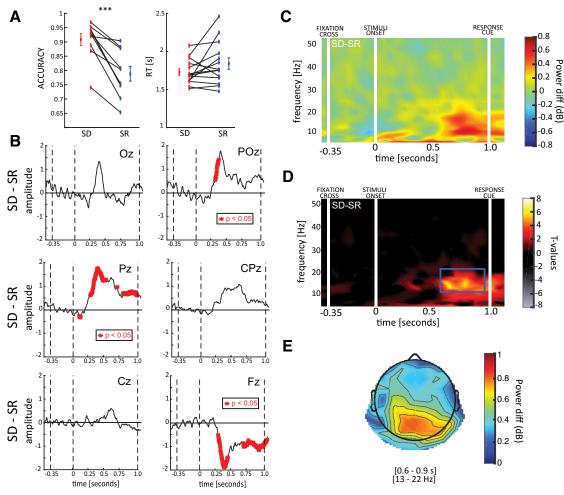


Figure 5. Pilot experiment results. **A**, Behavioral results of the pilot experiment: left and right panel show accuracy and RTs for SD (red) and SR (blue) tasks. Differently than in the main task, in the pilot experiment participants performed significantly better in the SD than in the SR task (compare the accuracy between **A** and Fig. 2B). ***indicates Bayes Factor > 100. **B**, Difference between SD and SR EPs. Red asterisk indicates time window significantly different from zero. **C**, Difference between SD and SR power spectra: white lines indicate the stimulus onset and the response cue. **D**, Testing the SD-SR difference against zero reveals a significant region in the low β band (13–21 Hz), before the response cue, in agreement with the results of the main experiment (Fig. 4). We reported a large effect size for this effect (one sample t test against zero averaging per each participant the values within the significant region, $t_{(13)} = 7.049$, ρ < 0.001, Cohen's d = 1.820). **E**, As in the main experiment, the SD-SR difference mostly involves occipital-parietal regions.

processes (Treisman and Gelade, 1980; Wolfe et al., 1989), and how increased activity in specific oscillatory components [i.e., low (22–34 Hz) and high (36–56 Hz) γ bands] are characteristic of these processes (Buschman and Miller, 2007; Phillips and Takeda, 2009). Accordingly, state-of-the-art computational models performing visual search and related tasks (e.g., instance segmentation) also employ attentional or recurrent mechanisms (Linsley et al., 2020), supporting the hypothesis that convolutional feedforward networks can benefit from recurrent mechanisms in solving visual reasoning tasks (Kreiman and Serre, 2020).

Computational evidence for the hypothesis that the SD task requires additional computational mechanisms beyond those needed to solve the SR task is provided by the results of the Siamese network simulations (Bromley et al., 1993). This feedforward network processes each stimulus item in a separate (CNN) channel and then

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



Previous work has shown that modulations of β band oscillations can be related to selective attention mechanisms (Buschman and Miller, 2007; Benchenane et al., 2011; Lee et al., 2013; Richter et al., 2018). Different attentional mechanisms may indeed be involved in the two tasks: the SR task could be solved by first grouping items and then determining the orientation of the group (Franconeri et al., 2012), whereas the SD task requires the individuation of the two items before comparison. In addition, our results are also consistent with differences in memory processes between the two tasks (de Fockert et al., 2001). One common assumption is that items that are grouped together (as in the SR task) occupy only one working memory slot (Franconeri et al., 2013; Clevenger and Hummel, 2014), whereas non-grouped items would each hold one slot, resulting in a larger working memory load. Previous literature showed that working memory can also be characterized by neuronal oscillatory signatures. Recent studies, for example, have demonstrated an interplay between β and γ band frequencies during working memory tasks (Lundqvist et al., 2016, 2018). Similarly, α and low β bands, not only increase with working memory load (Pesonen et al., 2007; Babiloni et al., 2004), but also in conjunction with the inhibition of competing visual memories in selective memory retrieval (Park et al., 2010; Waldhauser et al., 2012). Besides, previous studies have reported that increased oscillatory activity in the α band is a signature of attentional processes, and it can predict the likelihood of successful trials in many tasks (Händel et al., 2011; Klimesch, 2012; Nelli et al., 2017); however, in our current study we did not investigate differences between correct and incorrect trials, but between different types of tasks (involving SR or sameness judgment), after controlling for task difficulty . This could explain why α band amplitude differences were less prominent in our study. All considered, several lines of evidence point toward β oscillations as crucially involved in both attention and working memory related processes. These processes, therefore, might be part of the additional computational mechanisms required for SD tasks compared with SR tasks. Future work could more directly compare the attention and memory dependence of each task in human subjects.

That feedforward neural networks are limited in their ability to solve simple visual reasoning tasks is already being recognized by the computer vision and neuroscience communities (Yamins et al., 2013; Rajalingham et al., 2018; Kar et al., 2019). Current CNN extensions include modules for integrating local and global features (Chen et al., 2018) as well as recurrent neural architectures (Yang et al., 2018). Our results suggest that the human visual system also deploys additional computations beyond feedforward processes to successfully solve visual reasoning tasks. Rhythmic cortical oscillations in the β band represent the signatures of these additional computations, which may involve selective attention and working memory.

References

Babiloni C, Babiloni F, Carducci F, Cappa SF, Cincotti F, Del Percio C, Miniussi C, Vito Moretti D, Rossi S, Sosta K, Rossini PM (2004) Human cortical rhythms during visual delayed choice reaction time

- tasks: a high-resolution EEG study on normal aging. Behavioural Brain Res 153:261–271.
- Benchenane K, Tiesinga PH, Battaglia FP (2011) Oscillations in the prefrontal cortex: a gateway to memory and attention. Curr Opin Neurobiol 21:475–485.
- Bernardo JM, Smith AFM (2009) Bayesian theory. Sussex, UK:John Wiley & Sons.
- Brainard DH (1997) The psychophysics toolbox. Spat Vis 10:433–436
- Bromley J, Bentz JW, Bottou L, Guyon I, Lecun , Moore C, Sackinger E, Shah R (1993) Signature Verification using a "Siamese" time delay neural network". Intern J Pattern Recognit Artif Intell 7:669-688
- Buschman TJ, Miller EK (2007) Top-down versus bottom-up control of attention in the prefrontal and posterior parietal cortices. Science 315:1860–1864.
- Chen X, Li LJ, Fei-Fei L, Gupta A (2018) Iterative visual reasoning beyond convolutions. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp 7239–7248.
- Clevenger PE, Hummel JE (2014) Working memory for relations among objects. Atten Percept Psychophys 76:1933–1953.
- Daniel TA, Wright AA, Katz JS (2015) Abstract-concept learning of difference in pigeons. Anim Cogn 18:831–837.
- de Fockert JW, Rees G, Frith CD, Lavie N (2001) The role of working memory in visual selective attention. Science 291:1803–1806.
- Delorme A, Makeig S (2004) EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. J Neurosci Methods 134:9–21.
- Fabiani M, Karis D, Donchin E (1986) P300 and recall in an incidental memory paradigm. Psychophysiology 23:298–308.
- Franconeri SL, Scimeca JM, Roth JC, Helseth SA, Kahn LE (2012) Flexible visual processing of spatial relationships. Cognition 122:210–227.
- Franconeri SL, Alvarez GA, Cavanagh P (2013) Flexible cognitive resources: competitive content maps for attention and memory. Trends Cogn Sci 17:134–141.
- GoogleResearch (2016) TensorFlow: a system for large-scale machine learning, Proceedings of the 12th USENIX conference on Operating Systems Design and Implementation, pp. 265–283.
- Händel BF, Haarmeier T, Jensen O (2011) Alpha oscillations correlate with the successful inhibition of unattended stimuli. J Cogn Neurosci 23:2494–2502.
- He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, December 2016, pp 770–778, Las Vegas (USA).
- Hochberg B (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc Series B Stat Methodol 57:289–300.
- Hollard VD, Delius JD (1982) Rotational invariance in visual pattern recognition by pigeons and humans. Science 218:804–806.
- Itthipuripat S, Cha K, Byers A, Serences JT (2017) Two different mechanisms support selective attention at different phases of training. PLoS Biol 15:e2001724.
- Itthipuripat S, Cha K, Deering S, Salazar AM, Serences JT (2018) Having more choices changes how human observers weight stable sensory evidence. J Neurosci 38:8635–8649.
- JASP Team (2018) JASP (version 0.8.6.0). [Computer Software]. Retrieved from http://jasp-stats.org.
- Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ (2019) Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior. Nat Neurosci 22:974–983.
- Kim J, Ricci M, Serre T (2018) Not-so-CLEVR: learning same-different relations strains feedforward neural networks. Interface Focus 8:20180011.
- Kingma DP, Ba JL (2015) Adam: a method for stochastic gradient descent. International Conference on Learning Representations, San Diego (USA).
- Klimesch W (2012) Alpha-band oscillations, attention, and controlled access to stored information. Trends Cogn Sci 16:606–617.



- Kok A (2001) On the utility of P3 amplitude as a measure of processing capacity. Psychophysiology 38:557–577.
- Kreiman G, Serre T (2020) Beyond the feedforward sweep: feedback computations in the visual cortex. Ann NY Acad Sci 1464:222–241.
- Krusemark EA, Kiehl KA, Newman JP (2016) Endogenous attention modulates early selective attention in psychopathy: an ERP investigation. Cogn Affect Behav Neurosci 16:779–788.
- Lee JH, Whittington MA, Kopell NJ (2013) Top-down beta rhythms support selective attention via interlaminar interaction: a model. PLoS Comput Biol 9:e1003164.
- Linsley D, Ashok AK, Govindarajan LN, Liu R, Serre T (2020) Stable and expressive recurrent vision models. Retrieved from http://arxiv.org/abs/2005.11362.
- Love J, Selker R, Verhagen J, Marsman M, Gronau QF, Jamil T, Rouder JN (2015) Software to sharpen your stats. APS Observer 28:27–29.
- Lundqvist M, Rose J, Herman P, Brincat SLL, Buschman TJJ, Miller EKK (2016) Gamma and beta bursts underlie working memory. Neuron 90:152–164.
- Lundqvist M, Herman P, Warden MR, Brincat SL, Miller EK (2018) Gamma and beta bursts during working memory readout suggest roles in its volitional control. Nat Commun 9:394.
- Maris E, Oostenveld R (2007) Nonparametric statistical testing of EEG- and MEG-data. J Neurosci Methods 164:177–190.
- Masson MEJ (2011) A tutorial on a practical Bayesian alternative to null-hypothesis significance testing. Behav Res Methods 43:679– 690.
- McEvoy LK, Smith ME, Gevins A (1998) Dynamic cortical networks of verbal and spatial working memory: effects of memory load and task practice. Cereb Cortex 8:563–574.
- Nash AJ, Fernandez M (1996) P300 and allocation of attention in dual-tasks. Int J Psychophysiol 23:171–180.
- Nelli S, Itthipuripat S, Srinivasan R, Serences JT (2017) Fluctuations in instantaneous frequency predict alpha amplitude during visual perception. Nat Commun 8:1–12.
- Park HD, Min BK, Lee KM (2010) EEG oscillations reflect visual short-term memory processes for the change detection in human faces. Neuroimage 53:629–637.
- Pesonen M, Hämäläinen H, Krause CM (2007) Brain oscillatory 4-30 Hz responses during a visual n-back memory task with varying memory load. Brain Res 1138:171–177.
- Phillips PJ, Yates AN, Hu Y, Hahn CA, Noyes E, Jackson K, Cavazos JG, Jeckeln G, Ranjan R, Sankaranarayanan S, Chen JC, Castillo CD, Chellappa R, White D, O'Toole AJ (2018) Face recognition

- accuracy of forensic examiners, superrecognizers, and face recognition algorithms. Proc Natl Acad Sci USA 115:6171–6176.
- Phillips S, Takeda Y (2009) Greater frontal-parietal synchrony at low gamma-band frequencies for inefficient than efficient visual search in human EEG. Int J Psychophysiol 73:350–354.
- Rajalingham R, Issa EB, Bashivan P, Kar K, Schmidt K, DiCarlo JJ (2018) Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks. J Neurosci 38:7255–7269.
- Richter CG, Coppola R, Bressler SL (2018) Top-down beta oscillatory signaling conveys behavioral context in early visual cortex. Sci Rep 8:6991.
- Serre T (2019) Deep learning: the good, the bad, and the ugly. Annu Rev Vis Sci 5:399–426.
- Stabinger S, Rodríguez-Sánchez A, Piater J (2016) 25 years of cnns: Can we compare to human abstraction capabilities? In International Conference on Artificial Neural Networks, pp 380–387, Springer Cham.
- Treisman A, Gelade G (1980) A feature integration theory of attention. Cogn Psychol 12:97–136.
- VanRullen R, Thorpe SJ (2001) Is it a bird? Is it a plane? Ultra-rapid visual categorisation of natural and artifactual objects. Perception 30:655–668.
- Van Voorhis S, Hillyard SA (1977) Visual evoked potentials and selective attention to points in space. Percept Psychophys 22:54–62.
- Vogels R (1999) Categorization of complex visual images by rhesus monkeys. Part 1: behavioural study. Eur J Neurosci 11:1223–1238.
- Waldhauser GT, Johansson M, Hanslmayr S (2012) α/β oscillations indicate inhibition of interfering visual memories. J Neurosci 32:1953–1961.
- Wasserman EA, Castro L, Freeman JH (2012) Same-different categorization in rats. Learn Mem 19:142–145.
- Watson AB, Pelli DG (1983) Quest: a Bayesian adaptive psychometric method. Percept Psychophys 33:113–120.
- Wolfe JM, Cave KR, Franzel SL (1989) Guided search: an alternative to the feature integration model for visual search. J Exp Psychol Hum Percept Perform 15:419–433.
- Yamins D, Hong H, Cadieu C, Dicarlo JJ (2013) Hierarchical modular optimization of convolutional networks achieves representations similar to macaque IT and human ventral stream. In Advances in neural information processing systems, pp 3093–3101, Lake Tahoe, Nevada, United States.
- Yang GR, Ganichev I, Wang XJ, Shlens J, Sussillo D (2018) A dataset and architecture for visual reasoning with a working memory. In: European Conference on Computer Vision, pp 729–745, Springer, Cham. Conference location: Munich, Germany.