



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/240520/>

Version: Accepted Version

---

**Article:**

Strzelczyk, D., Clayson, P.E., Sigurdardottir, H.M. et al. (2026) Contralateral delay activity as a marker of visual working memory capacity: a multi-site registered replication. *Cortex*. ISSN: 0010-9452

<https://doi.org/10.1016/j.cortex.2026.04.006>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Journal Pre-proof



Contralateral delay activity as a marker of visual working memory capacity: a multi-site registered replication

Dawid Strzelczyk, Peter E. Clayson, Heida Maria Sigurdardottir, Faisal Mushtaq, Yuri G. Pavlov, H  l  ne Devillez, Anton Lukashevich, Harold A. Rocha, Yong Hoon Chung, Kevin M. Ortego, Viola S. St  rmer, Jos   C. Garc  a Alanis, Christoph L  ffler, Anna-Lena Schubert, Anna Lena Biel, Samuel A. Birkholz, Emily M. Johnson, Jeffrey S. Johnson, Zitong Lu, Yong Min Choi, Eva Lout, Julie D. Golomb, Shuangke Jiang, Myles Jones, Eda Mizrak, Claudia C. von Bastian, Niko A. Busch, Charline Peylo, Larissa Behnke, Yannik Hilla, Maro G. Machizawa, William X.Q. Ngiam, Edward K. Vogel, Nicolas Langer

PII: S0010-9452(26)00113-9

DOI: <https://doi.org/10.1016/j.cortex.2026.04.006>

Reference: CORTEX 4340

To appear in: *Cortex*

Received Date: 24 February 2026

Revised Date: 8 April 2026

Accepted Date: 8 April 2026

Please cite this article as: Strzelczyk D, Clayson PE, Sigurdardottir HM, Mushtaq F, Pavlov YG, Devillez H, Lukashevich A, Rocha HA, Chung YH, Ortego KM, St  rmer VS, Garc  a Alanis JC, L  ffler C, Schubert A-L, Biel AL, Birkholz SA, Johnson EM, Johnson JS, Lu Z, Choi YM, Lout E, Golomb JD, Jiang S, Jones M, Mizrak E, von Bastian CC, Busch NA, Peylo C, Behnke L, Hilla Y, Machizawa MG, Ngiam WXQ, Vogel EK, Langer N, Contralateral delay activity as a marker of visual working memory capacity: a multi-site registered replication, *Cortex*, <https://doi.org/10.1016/j.cortex.2026.04.006>.

This is a PDF of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability. This version will undergo additional copyediting, typesetting and review before it is published in its final form. As such, this version is no longer the Accepted Manuscript, but it is not yet the definitive Version of Record; we are providing this early version to give early visibility of the article. Please note that Elsevier's sharing policy for the Published Journal Article applies to this version, see: <https://www.elsevier.com/about/policies-and->

[standards/sharing#4-published-journal-article](#). Please also note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2026 The Author(s). Published by Elsevier Ltd.

# Contralateral delay activity as a marker of visual working memory capacity: a multi-site registered replication

Dawid Strzelczyk<sup>1,2</sup>, Peter E. Clayson<sup>3</sup>, Heida Maria Sigurdardottir<sup>4</sup>, Faisal Mushtaq<sup>5</sup>, Yuri G. Pavlov<sup>6</sup>, H el ene Devillez<sup>4</sup>, Anton Lukashovich<sup>4</sup>, Harold A. Rocha<sup>3</sup>, Yong Hoon Chung<sup>7</sup>, Kevin M. Ortego<sup>7</sup>, Viola S. St ormer<sup>7</sup>, Jos e C. Garc ia Alanis<sup>8</sup>, Christoph L offler<sup>8</sup>, Anna-Lena Schubert<sup>8</sup>, Anna Lena Biel<sup>9</sup>, Samuel A. Birkholz<sup>10</sup>, Emily M. Johnson<sup>10</sup>, Jeffrey S. Johnson<sup>10</sup>, Zitong Lu<sup>11</sup>, Yong Min Choi<sup>11</sup>, Eva Lout<sup>11</sup>, Julie D. Golomb<sup>11</sup>, Shuangke Jiang<sup>12,20</sup>, Myles Jones<sup>12</sup>, Eda Mizrak<sup>12,14</sup>, Claudia C. von Bastian<sup>12</sup>, Niko A. Busch<sup>9</sup>, Charline Peylo<sup>13</sup>, Larissa Behnke<sup>13</sup>, Yannik Hilla<sup>13</sup>, Maro G. Machizawa<sup>15</sup>, William X. Q. Ngiam<sup>16,17</sup>, Edward K. Vogel<sup>16,17</sup>, Nicolas Langer<sup>1,2,18\*</sup>

<sup>1</sup> Methods of Plasticity Research, Department of Psychology, University of Zurich, Zurich, Switzerland

<sup>2</sup> Neuroscience Center Zurich (ZNZ), University of Zurich and ETH Zurich, Zurich, Switzerland

<sup>3</sup> University of South Florida, Tampa, Florida, USA

<sup>4</sup> University of Iceland, Reykjavik, Iceland

<sup>5</sup> University of Leeds, Leeds, UK

<sup>6</sup> University of Tuebingen, Tuebingen, Germany

<sup>7</sup> Dartmouth College, Hanover, New Hampshire, USA

<sup>8</sup> University of Mainz, Germany

<sup>9</sup> Institute of Psychology, University of M unster, Germany

<sup>10</sup> Department of Psychology, North Dakota State University, Fargo, ND, USA

<sup>11</sup> The Ohio State University, Columbus, Ohio, USA

<sup>12</sup> University of Sheffield, Sheffield, UK

<sup>13</sup> Neuropsychology and Cognitive Neuroscience, Department of Psychology, University of Zurich, Zurich, Switzerland

<sup>14</sup> University of Oxford, Oxford, UK

<sup>15</sup> Digital Brain Science Laboratory, Xiberlinc Inc., Tokyo, Japan

<sup>16</sup> Department of Psychology, University of Chicago, Illinois, USA

<sup>17</sup> Institute of Mind and Biology, University of Chicago, Illinois, USA

<sup>18</sup> Center of Reproducible Science, University of Zurich, Zurich, Switzerland

<sup>19</sup> Cognitive Psychology, Department of Psychology, University of Zurich, Zurich, Switzerland

\* n.langer@psychologie.uzh.ch

## 1 Abstract

2

3 The contralateral delay activity (CDA) is a widely used electrophysiological marker of visual  
4 working memory (VWM), yet recent work has questioned whether typical sample sizes in CDA  
5 studies are sufficient to robustly detect set size effects and brain-behavior correlations. As part  
6 of the #EEGManyLabs initiative, the present multi-site replication study aimed to rigorously  
7 test replicability of the key findings of Vogel and Machizawa (2004) using a large sample of  
8 304 participants across 10 laboratories and a preregistered analysis plan. We replicated the  
9 expected contralateral-ipsilateral asymmetry and observed increases in CDA amplitude from  
10 set size 2 to 4 and from set size 2 to 6. In contrast, the hypothesized positive correlation  
11 between the CDA increase from set size 2 to 4 and individual VWM capacity was not replicated  
12 in the preregistered meta-analytic correlation. Across different pipelines and statistical  
13 analyses, the meta-analytic correlation estimate was small ( $r = 0.15$ ) and substantially  
14 attenuated relative to the original effect size in Vogel and Machizawa (2004) study ( $r = 0.78$ ).  
15 To contextualize these findings, we applied a funnel-plot diagnostic combining published  
16 effects with the #EEGManyLabs data, indicating small-study inflation and publication bias.  
17 Taken together, our results indicate that reports of strong correlations between CDA amplitude  
18 and VWM capacity may have been overestimated, in part because statistically significant  
19 findings were selectively reported. Our results highlight the importance of open science  
20 practices, including well-powered, preregistered studies with transparent data and analysis  
21 pipelines, in order to characterize the magnitude and robustness of individual-difference  
22 associations in psychophysiology.

23

24

25

26

27

28

## 1 Introduction

2

3 Visual working memory (VWM) is a temporary storage system that holds information that can  
4 be accessed and manipulated by higher cognitive functions (Luck & Vogel, 2013). Visual  
5 working memory is considered a central construct in cognitive neuroscience and is a putative  
6 intermediary for information transfer (Atkinson & Shiffrin, 1968; A. Baddeley, 2003; A. D.  
7 Baddeley, 1986; A. D. Baddeley & Hitch, 1974; Cowan et al., 2005; Cowan & Morey, 2006;  
8 Liesefeld & Müller, 2019; Olivers, 2008), thereby facilitating various cognitive functions  
9 including reading comprehension (Caplan & Waters, 1999; Daneman & Carpenter, 1980; Lotfi  
10 et al., 2022; Martin & Romani, 1994; Wang et al., 2022), planning and problem-solving (Cowan  
11 et al., 2005; Miyake & Shah, 1999; Naveh-Benjamin & Cowan, 2023), and learning new skills  
12 (A. Baddeley et al., 1988; A. D. Baddeley et al., 2017; Cowan, 2014; Gathercole & Baddeley,  
13 1989; Jongbloed-Pereboom et al., 2019; von Bastian et al., 2022). Researchers have been  
14 using electroencephalography (EEG) to understand the neural correlates of VWM in real-time.  
15 A commonly used EEG measure of VWM is an event-related potential (ERP) called the  
16 contralateral delay activity (CDA). This EEG signal has also been referred to in other studies  
17 as Contralateral Negative Slow Wave (CNSW) by Klaver et al. (1999), Sustained Posterior  
18 Contralateral Negativity (SPCN) by Brisson and Jolicoeur (2007), Perron et al. (2009), and  
19 Contralateral Search Activity (CSA) by Emrich et al. (2009). These different terms all refer to  
20 the same visual working memory correlate. Hence, we will maintain the use of the term CDA  
21 throughout the remainder of the paper.

22

23 The CDA is a difference wave constructed by subtracting ipsilateral from contralateral activity  
24 related to the to-be-remembered items. Since items in studies analyzing experimental effects  
25 on the CDA are generally shown bilaterally, while only those on one side of the screen are  
26 supposed to be memorized, the idea of the subtraction is to eliminate any activity related to  
27 early perceptual and low-level processing by assuming that they equally affect ipsilateral and  
28 contralateral ERPs. Activity over the contralateral hemisphere tends to be more negative than  
29 ipsilateral activity during VWM retention (Luria et al., 2016; Ngiam et al., 2021; Vogel &  
30 Machizawa, 2004). Thus, it has been suggested that the CDA reflects the neural activity  
31 related to the maintenance of information in VWM, and studies have shown that the amplitude  
32 and duration of the CDA are linked to the amount of information stored in working memory.

33

34 In a seminal paper from Vogel and Machizawa (2004), the authors demonstrated that the CDA  
35 amplitude increases with the number of items stored in VWM and plateaus at around 3 to 4  
36 items, consistent with the typical adult working memory capacity (Forsberg et al., 2023). More

1 importantly, this study showed that the increase in the CDA amplitude with greater memory  
2 load correlated with individual VWM performance (Vogel & Machizawa, 2004). Specifically,  
3 individuals with high VWM capacity exhibited a larger increase in the CDA when attempting to  
4 memorize 4 compared to 2 items. In this study, the CDA was elicited using a color change  
5 detection task (Vogel & Machizawa, 2004). The task involves presenting participants with a  
6 central arrow cue that indicates whether participants need to memorize items on the left or  
7 right of the screen center. The cue is followed by a bilateral stimulus array with equal numbers  
8 of colored squares shown on each side (set size 1 to 10). After a short retention phase,  
9 participants are presented with a second array and asked to indicate whether any of the  
10 squares on the cued side changed color (Figure 2). The lateralized color change detection  
11 task is now a widely used paradigm to examine visual working memory processes (Feldmann-  
12 Wüstefeld, 2021; Luria et al., 2016) and has been explored in several variations such as  
13 different set sizes, including distractions, retro-cueing and using different shapes and colors  
14 (Feldmann-Wüstefeld, 2021; Feuerstahler et al., 2019; Roy & Faubert, 2023; Schneider et al.,  
15 2017).

16  
17 The finding that the CDA amplitude is sensitive to how much visual information is to be  
18 remembered has been replicated in numerous studies (Asp et al., 2021; Brady et al., 2016;  
19 Hakim et al., 2019; Heuer & Schubö, 2016; Quirk et al., 2020; Unsworth et al., 2015).  
20 Furthermore, several studies have validated the positive correlation between the CDA  
21 amplitude increase and VWM capacity (Adam et al., 2018; Feldmann-Wüstefeld, 2021;  
22 Villena-González et al., 2020). In the review paper of Luria et al. (2016), the authors conducted  
23 a meta-analysis from 11 previous studies and reported an aggregated correlation of  $r = 0.596$ .  
24 However, a recent study indicated that the typical numbers of subjects and trials for CDA  
25 experiments seen in the literature may be underpowered for detecting set size differences  
26 (Ngiam et al., 2021).

27  
28 The insufficient power issue is even more pressing for the correlation between the VWM  
29 capacity and the CDA amplitude increase. Critically, Schönbrodt and Perugini (2013)  
30 demonstrated that correlation estimates typically stabilize at a sample size of approximately  
31 250 subjects (Schönbrodt & Perugini, 2013). Except for one large study ( $N = 171$ ; Unsworth  
32 et al., 2015), the average sample size of previous studies investigating the relationship  
33 between VWM capacity and the CDA amplitude was 32 subjects (range 12-83 subjects for 12  
34 studies; Luria et al, 2016). Finally, the inherent flexibility in EEG analysis, including analysis of  
35 the CDA, leaves many decisions up to the researcher. This leaves open the possibility to  
36 exploit these researchers' degrees of freedom (i.e., the garden of forking paths; Gelman &  
37 Loken, 2013), either intentionally or unintentionally. Such practices can lead to erroneous

1 inferences and perpetuate replication problems in cognitive neuroscience (Clayson et al.,  
2 2019; Luck & Gaspelin, 2017).

3

4 To address this issue, the #EEGManyLabs project was initiated (Pavlov et al., 2021). The  
5 #EEGManyLabs initiative highlights the importance of replication in science and the need for  
6 rigorous research methods to increase confidence in prominent effects. The #EEGManyLabs  
7 project aims to replicate pivotal EEG studies which had a critical impact on the cognitive and  
8 affective neuroscience community. Importantly, the #EEGManyLabs project is designed to  
9 address some of the limitations of previous replication efforts by using a large sample of  
10 participants, standardized procedures, and a pre-registered analysis plan (i.e., Registered  
11 Report; Pavlov et al., 2021).

12

13 As part of the #EEGManyLabs project, the current study aimed to contribute to the existing  
14 literature on VWM and the CDA by conducting a robust multi-site, large-scale replication of  
15 Vogel and Machizawa's (2004) seminal study. The present study was chosen for replication  
16 by a global consortium of EEG specialists owing to its scientific significance (for further  
17 information on the selection process, refer to Pavlov et al., 2021). In accordance with the  
18 #EEGManyLabs project, this Registered Report closely adhered to the original study design  
19 and ensured adequate statistical power with a large sample size. The present study also  
20 followed preregistered analysis steps to ensure the integrity of the direct replication and  
21 statistical inferences (Paul et al., 2021). To this end, Experiment 3 of the original study was  
22 replicated using three set sizes (2, 4, and 6 items per side) to examine whether CDA amplitude  
23 varies as a function of memory load. In line with the original study, the following hypotheses  
24 were tested:

25

26 [H1.1] The CDA amplitude increases from arrays of 2 items per side to arrays of 4 items per  
27 side.

28 [H1.2] The CDA amplitude increases from arrays of 2 items per side to arrays of 6 items per  
29 side.

30 [H1.3] The CDA amplitude for 4 items and 6 items is equivalent.

31

32 Additionally, the study examined whether the CDA amplitude is related to performance on the  
33 change detection task:

34

35 [H2.1] Subjects' VWM capacity (measured behaviorally) is positively correlated with the CDA  
36 amplitude increase from 2 to 4 items.

1 [H2.2] Subjects' VWM capacity (measured behaviorally) is not correlated with the CDA  
2 amplitude increase from 4 to 6 items.

3

4 Finally, replication success was evaluated separately for each hypothesis. For hypotheses  
5 predicting directional effects, replication success was defined as a statistically significant  
6 random-effects meta-analytic estimate in the same direction as in the original study, combining  
7 results across laboratories. For hypotheses predicting equivalence or the absence of an  
8 association, replication success was evaluated using equivalence testing.

Journal Pre-proof

## 2. Methods

The protocol for this replication was developed in consultation with the original authors (co-authors of the present work, EV, MM). The current document is a Stage 2 Registered Report that follows guidelines for open science in psychophysiological research as outlined by Garrett-Ruffin et al. (2021). The Stage-1 Registered Report outlining the preprocessing and analysis steps is available at: <https://doi.org/10.31234/osf.io/shdea> (Strzelczyk et al., 2023). All raw EEG and behavioral datasets, after marker harmonization, anonymization, and including full datasets that were later excluded from analysis, are openly accessible at: [https://gin.g-node.org/EEGManyLabs/EEGManyLabs\\_Replication\\_VogelMachizawa2004\\_Raw](https://gin.g-node.org/EEGManyLabs/EEGManyLabs_Replication_VogelMachizawa2004_Raw) and [https://gin.g-node.org/EEGManyLabs/EEGManyLabs\\_Replication\\_VogelMachizawa2004\\_Processed](https://gin.g-node.org/EEGManyLabs/EEGManyLabs_Replication_VogelMachizawa2004_Processed). All preprocessing and analysis scripts used in the present report are available on GitHub: <https://github.com/ksgfan/EEGManyLabs>. Each site has obtained approval from the local ethics committee to conduct the study and share data.

The institution abbreviations used throughout the manuscript are: DART (Dartmouth College, USA), USF (University of South Florida, USA), JGU (Johannes Gutenberg University Mainz, Germany), WWU (University of Münster, Germany), NDSU (North Dakota State University, USA), OSU (The Ohio State University, USA), UI (University of Iceland, Iceland), TUOS (University of Sheffield, UK), UZH-MPR (University of Zurich, Methods of Plasticity Research, Switzerland), and UZH-NCN (University of Zurich, Neuropsychology and Cognitive Neuroscience, Switzerland).

### Known differences from the original study

Table 1. Details on original, replication and alternative pipelines. Deviations from the original preprocessing are highlighted in blue in the direct replication pipeline.

Offline Processing Step	Original and current study parameters for the direct preprocessing pipeline	#EEGManyLabs: advanced preprocessing pipeline
Offline Filter	(1) Hardware online filter: bandpass of 0.01-80 Hz (half-power cutoff, Butterworth filters) (2) 35 Hz LP only for plots	(1) Offline bandpass of 0.01-80 Hz (half-power cutoff, Butterworth filters) (2) 35 Hz LP for plotting
Line noise removal	no line noise removal <a href="#">ZapLine method</a> <sup>1</sup>	ZapLine method
Ocular artifact rejection	(1) Trials containing ocular artifacts were removed (i.e., blinks or eye	If eye-tracker recording is available, we excluded trials with eye movements larger

	movements larger than 1 degree). A heuristic for 1 visual degree was used (25 microvolt bipolar HEOG amplitude threshold; adjusting the threshold for each subject based on visual inspection). A calibration paradigm was used to estimate the subject specific amplitude representing 1 visual degree	than 1 degree. If no eye-tracker is available, we identified ocular artifacts using a tailored subject-specific amplitude threshold for the EOG electrodes, which was obtained from the saccadic calibration task.
	(2) Blinks: unipolar VEOG >? microvolt Blinks: unipolar VEOG >90 microvolt	
Artifact Rejection	(1) peak-to-peak amplitude >200 microvolt (2) Visual inspection was used to identify and exclude trials containing movement artifacts or blocking	(1) Peak-to-peak amplitude >200 microvolt (2) bad trial identification method introduced by Adam, Robison, and Vogel (2018)
Bad Channel Identification	(1) peak-to-peak amplitude >75 microvolt. Bad channels were not interpolated, but artifactual trials were rejected (2) Visual inspection	(1) Correlation below 0.85 with neighboring channels (2) 4 SD or more line noise relative to signal than all other channels (3) Blocking longer than 5 s
Bad Channel Interpolation	NA	Spherical spline interpolation
Reference	Algebraic average of the left and right mastoids	Algebraic average of the left and right mastoids
CDA time interval	300 - 900 ms after memory onset	300 - 900 ms after memory onset
Baseline Interval	-200 - 0 ms	-200 - 0 ms
Region of Interest	left electrode cluster: P3, T5/P7, O1. right electrode cluster: P4, T6/P8, O2	left electrode cluster: P3, T5/P7, O1. right electrode cluster: P4, T6/P8, O2
CDA time interval	retention phase (i.e., 300 - 900 ms after the onset of memory array)	retention phase (i.e., 300 - 900 ms after the onset of memory array)
Set Size	Experiment #3: 2,4,6	2, 4, 6
Visual Memory Capacity	K & d'	K & d'

1 Note. Deviations from the original study in the direct preprocessing pipeline are shown in blue.

2 <sup>1</sup> Several labs are recording the task with an eye-tracker, which induces line noise. Therefore, we  
3 decided to use ZapLine to reduce the line noise.

4

## 5 Sample size & Inclusion criteria

6

7 Participants were recruited from universities or nearby communities. The study only included  
8 individuals between 18 and 35 years free from any diagnosed psychiatric or neurological  
9 disorders and with intact color vision. We acquired demographics (i.e., age, gender),  
10 handedness (Edinburgh Handedness Inventory; Oldfield, 1971) and education level based on  
11 International Standard Classification of Education (ISCED;  
12 <http://uis.unesco.org/en/topic/international-standard-classification-education-isced>).

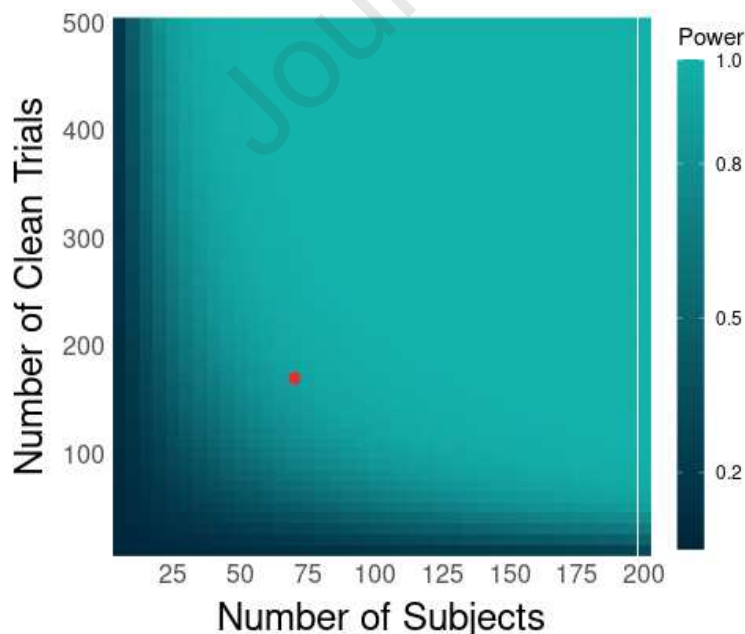
13

1 The required sample size was estimated for each hypothesis: For hypothesis #1, we used the  
 2 CDA power calculator (<https://williamngiam.shinyapps.io/CDAPower/>; Ngiam et al., 2021) to  
 3 estimate the required sample size to detect a set-size effect between set size 2 and 4 (which  
 4 is similar as between set size 2 and 6). With a minimum of 170 clean trials per condition (i.e.,  
 5 excluding subjects with a bad trial rate > 30%) and 90% power, the estimated number of  
 6 subjects required is 70 (see Figure 1).

7

8 The following procedure was conducted to estimate the required sample size to investigate  
 9 the correlation between VWM capacity and the CDA amplitude difference (i.e., hypothesis #2):  
 10 In the original study, 36 participants were recruited and the subjects' VWM capacity was  
 11 correlated with the CDA amplitude increase between 2 and 4 items with a correlation estimate  
 12 of  $r = .78$  (Vogel & Machizawa, 2004). The power analysis showed that with an alpha level of  
 13  $.02$  and an assumed effect size of 50% (i.e.,  $r = .39$ ) of the original study, a sample size of  $N$   
 14  $= 68$  is required to achieve 90% power in detecting the effect. For the sample size calculation  
 15 of hypothesis #2, we used the R package "pwr" (Champely, 2020) (`pwr.r.test(r = 0.39, sig.level`  
 16 `= 0.02, power = 0.9, alternative = "greater"`). However, according to Schönbrodt and Perugini  
 17 (2013), correlation starts to stabilize at the sample size  $N = 250$ . As the #EEGManyLabs is  
 18 open for any lab to participate, we decided that each participating lab (i.e.,  $N = 10$ ) should  
 19 recruit 25 participants, resulting in 250 participants in total, which provides sufficient power to  
 20 investigate both hypotheses.

21



22

23 Figure 1. CDA power calculation. We estimated the required sample size for the set size effect between  
 24 set size 2 and 4, assuming at least 170 trials (i.e., maximum of 30% bad trials) and 90% power. The  
 25 estimated number of subjects required is 70 (red dot).

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

## Exclusion criteria

The color change detection task requires the participants to discriminate between colored squares, therefore color blindness is a critical exclusion criterion. We tested the color-vision with an online color-vision test (<https://colormax.org/color-blind-test/>). Participants with scores below 11 out of 12 correct responses were considered to have a color-vision deficiency and were excluded from further analyses.

## Exclusion criteria for direct replication

Following the original study, we excluded trials with eye movements, blinks, and blocking (amplifier saturation after drift). To identify eye movements and blinks, horizontal electrooculography (EOG) was concurrently recorded. Contaminated trials were identified by large ( $>1^\circ$ ) eye movements (Vogel & Machizawa, 2004).

In the original study, the authors used a heuristic for  $1^\circ$  horizontal eye movements and a fixed amplitude threshold for each subject. In this replication study, we deviated from this procedure and calculated the  $1^\circ$  horizontal eye movement amplitude for each participant in order to more accurately estimate an individual's amplitude threshold. To determine the individual participant's exclusion amplitude threshold, which reflects  $1^\circ$  horizontal eye movement, there was a separate horizontal EOG saccade calibration task prior to the main experiment (developed by K.M.Ortego, co-author). This task involves participants making saccades to left and right targets on the screen. Participants started each trial by fixating on the center of the screen. Following a key press there was a jittered interval between 1200~1600 ms and a saccade target (a red disk;  $0.6^\circ$  in size) appeared either  $3^\circ$  or  $6^\circ$  away from the fixation on the left or right side of the screen along the horizontal midline. Participants were instructed to make a saccade to the target location as soon as it appeared and to press a space bar once they had successfully made the saccade. There were 15 trials per condition, resulting in 60 trials total. The data from the saccade calibration paradigm was preprocessed by (1) bandpass filtering the data from 0.1-40 Hz; (2) epoching from -200 to +600 ms with respect to the onset of the saccade target; and (3) baseline correcting using a pre-stimulus baseline interval of -200 to 0 ms. Given previous research showing the saccade onset latency being  $\sim 200$  ms (Westheimer, 1954a, 1954b), horizontal EOG (i.e.,  $\text{HEOG} = \text{HEOGR} - \text{HEOGL}$ ) channel amplitudes from horizontal saccades were averaged during the 300~400 ms interval across the left and right conditions. The  $1^\circ$  horizontal eye movement amplitude threshold was then

1 calculated by extrapolating from 3° and 6° eye movements (estimating the linear regression  
2 curve using fitype function in MATLAB) as previous reports have shown the HEOG amplitudes  
3 and the size of saccades have a consistently linear relationship (Luck, 2014). We did not  
4 measure the 1° eye movement directly, as a pilot study demonstrated that estimating the 1°  
5 eye movement is more error prone and has too much variability. Furthermore, blinks were  
6 detected by using an amplitude threshold (>50 microvolt) in the unipolar VEOG channel. In  
7 addition, a segment was marked as bad if any electrodes of interest (i.e., HEOG, P3, T5/P7,  
8 O1; P4, T6/P8, O2) had a peak-to-peak amplitude >200 microvolt within one segment. Finally,  
9 visual inspection was used to identify bad trials. For an overview of the exclusion criteria and  
10 analysis pipeline see the Table 1. If more than 30% of trials (all set-size conditions combined)  
11 had to be rejected by these combined criteria, the subjects were excluded from further analysis  
12 to assure sufficient number of trials (see sample size calculation).

13

#### 14 Alternative analysis pipelines

15

16 For the alternative pipelines (Table 1), to identify eye movements and blinks, all labs recorded  
17 horizontal and vertical EOG and several labs additionally recorded eye tracking data (see  
18 Table 2). Trials contaminated by eye movements larger than 1° (Vogel & Machizawa, 2004)  
19 were identified based on eye tracking data if available (i.e., trials containing eye movements  
20 larger than 1°), and otherwise based on EOG data as described in the previous section. In  
21 addition to the bad trial identification methods described in the direct replication, we also  
22 utilized the bad trial identification method introduced by Adam et al. (2018; see below). Again,  
23 if more than 30% of trials for a specific set size had to be rejected by these combined criteria,  
24 the subjects were excluded from further analysis to assure sufficient number of trials (see  
25 sample size calculation).

26

#### 27 Procedure

28

29 Upon their arrival, participants received a brief overview of the experiment and were asked to  
30 give their informed written consent for participating in the study and allowing data sharing.  
31 Next, the participants were asked to fill out a short questionnaire regarding their history of  
32 psychiatric and neurological disorders, handedness and educational level, and carry out an  
33 online color-blind test (<https://colormax.org/color-blind-test/>). Subsequently, the participants  
34 were comfortably seated in a chair. If available, the experiment was conducted in a sound-  
35 and electrically shielded Faraday recording cage. Some cages were equipped with a chinrest  
36 to minimize head movements. A cap with electrodes was placed on the participant's head and

1 impedances were checked if provided by the EEG amplifier system and improved if necessary  
2 (see Table 2 for details). As this project is part of a wider initiative on replicability in EEG  
3 (#EEGManyLabs), several of the laboratories in this replication also collected resting state  
4 EEG data together with some personality measures (<https://osf.io/sp3ck/>, Pavlov et al., 2021).  
5 Neither resting EEG nor personality data were analyzed in the current study but were merged  
6 across sites as part of a future replication project to be reported elsewhere. Participants first  
7 completed an EOG saccade calibration task, after which the color change detection task  
8 began. The expected duration of the entire experiment was approximately 120 minutes. Upon  
9 completion of the examination, participants received compensation or credit for their  
10 participation.

11

## 12 Experimental Paradigm

13

14 The color change detection task was identical to the task used in the original study (Vogel &  
15 Machizawa, 2004). The paradigm was implemented in MATLAB, using the PsychToolbox  
16 extensions (Brainard, 1997; Pelli, 1997). Each trial of the task started with a blank screen  
17 presented for 1500 ms. Then, a central arrow appeared for 200 ms, indicating which side of  
18 the screen the participant should pay attention to. This was followed by another fixation period  
19 of a random time interval between 300 ms and 400 ms. Afterwards, a memory set was  
20 presented for 100 ms, which consisted of either 2, 4, or 6 colored squares on each side of the  
21 screen. Participants were instructed to only memorize the part of the memory set indicated by  
22 the arrow. This was followed by a 900 ms retention interval with a blank screen and a fixation  
23 cross (see Figure 2). Finally, a test array was presented for 2000 ms, and the participants  
24 were asked to indicate whether the test array was identical to the previous memory array ("no-  
25 change" trial) or whether the test array was different by one color ("change" trial).

26

27 The participants indicated whether a change occurred by pressing either the A or L button on  
28 a keyboard. The button they pressed depended on the instruction they were given, with half  
29 of the participants being instructed to press the A button for a change and the L button for no  
30 change, while the other half was instructed to do the opposite. Additionally, the participants  
31 were instructed to use their left hand to press the A button and their right hand to press the L  
32 button. During the task, participants were asked to focus their gaze on the fixation cross in the  
33 center of the screen until the probe appeared.

34

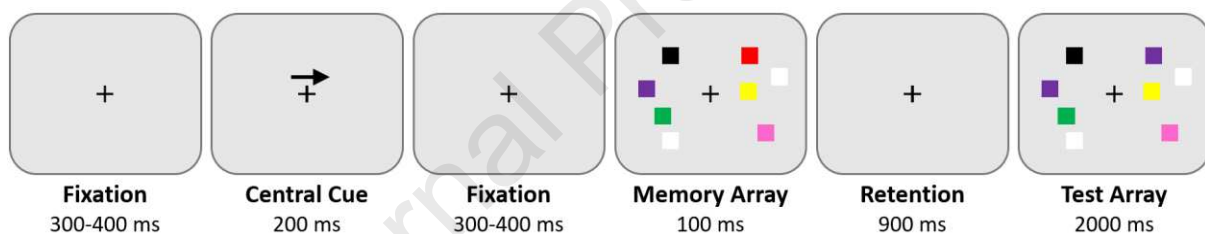
35 All stimuli were displayed within 2 regions that were  $4^\circ \times 7.3^\circ$  in size and were located  $3^\circ$  to  
36 the left and right of a central fixation cross on a gray background ( $8.2 \text{ cd m}^{-2}$ ). Each memory

1 array consisted of 2, 4 or 6 colored squares ( $0.65^\circ \times 0.65^\circ$ ) in each visual field. The squares  
 2 were chosen at random from a set of seven highly distinct colors (red, blue, violet, green,  
 3 yellow, black and white), and a specific color appeared no more than twice in a single array.  
 4 In other words, a specific color could be displayed in both hemifields but never twice within a  
 5 hemifield. The positions of the stimuli were randomized on each trial, with the restriction that  
 6 the distance between squares within a visual field was at least  $2^\circ$  (center to center). In 50% of  
 7 the trials, the color of one square in the test array on the cued side was different from the  
 8 corresponding square in the memory array, while in the remaining trials, the colors of the 2  
 9 arrays were identical.

10

11 The task was divided into five blocks, each containing 144 trials (i.e., 720 trials per subject  
 12 and 240 trials per condition and subject). The cue direction (left or right) and set size (2, 4 or  
 13 6 items on each side of the screen) were randomly varied throughout the trials to ensure a  
 14 balanced distribution of all conditions in each block. In line with the original study, no training  
 15 exercise was conducted prior to the main task.

16



17

18 Figure 2. Lateralized color change detection task. The figure is illustrative and not to scale.

19

## 20 Neurophysiological Data Acquisition

21

22 The replicating labs used a range of EEG systems and, where available, eye trackers.  
 23 Acquisition details are provided in Table 2. All labs provided the raw data to Zurich's Lab (UZH-  
 24 MPR), where it was preprocessed and analyzed.

25

26 Table 2. Data acquisition settings for each laboratory.

Lab	Screen type; size; ratio; refresh rate	Stimulus presentation language	Distance between chinrest and monitor	EEG system; number of channels; sampling rate	Reference; grounding	Impedances	Eye tracker; sampling rate	HEOG	Faraday cage	Soundproof or sound attenuated recording room
Dartmouth College	VPIxx; 540x300 mm; 1090 x 1080; 120 Hz	Psychtoolbox 3.0.18	45 cm	BrainVision; 32 channels; 500 Hz	Right mastoid; Fpz	Kept below 10 KOhm	No	Yes	Yes	Yes
University of South Florida	Dell p2314h, 23" widescreen, 60 Hz	Psychtoolbox 3.0.18	65 cm	Magstim EGI, 128 channels, 500 Hz	Cz; PCz	Kept below 50 KOhm	No	Yes	No	No
University of	Eizo	Psychtoolbox	67 cm	BrainProducts	FCz; Fpz	Kept below 10	No	Yes	Yes	Yes

Mainz	ColorEdge CS2420; 24.1" diag; 1920x1200; 60 Hz	3.0.18		; 64 channels; 1000 Hz		KOhm					
University of Münster	Viewpixmap G 1920 x 1080 120Hz	Psychtoolbox 3.0.18	86	Biosemi; 64 + 3 chans; 1024Hz	Reference free; GND adjacent to POz	Not available with Biosemi	Eye Link, 500 Hz	Yes	no	Yes	
North Dakota State University	ASUS ROG Strix XG27AQ 27"; 2560 x 1440;	Psychtoolbox 3.0.18	50 cm	Biosemi; 64 + 8 chans; 512 Hz	Reference free; GND adjacent to POz	Not available with Biosemi	Eye Link 1000, 500 Hz	Yes	Yes	Yes	
The Ohio State University	BENQ XL2420-B; 1920 x 1080; 120 Hz	Psychtoolbox 3.0.18	80cm	Brain Vision; 32 channels; 1000Hz	Cz, Fpz	Kept below 20 KOhm	EyeLink 1000; 500 Hz	Yes	Yes	Yes	
Icelandic Vision Lab, University of Iceland	2560*1440 60 Hz ASUS PG278QR 27"	Psychtoolbox 3.0.18	57cm (nasion to screen distance; "no chinrest for pilot)	BrainVision; 32 channels; 1000 Hz	Fz, Fpz	Kept below 15 KOhm	No	Yes	No	Yes	
University of Sheffield	Iiyama G-master GB2488HS U; 531.4 x 298.9mm; 1920 x 1080; 144 Hz	Psychtoolbox 3.0.18	50cm (nasion to screen distance; "no chinrest for pilot)	Biosemi; 64 channel; Recorded at 2048 Hz and then downsampled to 512 Hz	Unreferenced/reference free; CMS/DRL adjacent to POz	Not Available (Only offset within $\pm 25$ mV)	No	Yes	Yes	Yes	
University of Zürich (UZH-MPR)	Philips 242E1; 540x414mm; 800x600; 100 Hz	Psychtoolbox 3.0.18	70 cm	ANT Neuro; 128 channels; 500 Hz	CPz; GND adjacent to M1	Kept below 20 KOhm	EyeLink 1000; 500 Hz	Yes	Yes	Yes	
University of Zürich (UZH-NCN)	HP Omen 27q; 2560 x 1440; 144 Hz	Psychtoolbox 3.0.18	70 cm; no chin rest	BrainProducts; 64 channels; 1000 Hz	FCz; Fpz	Kept below 15 KOhm; ground and reference electrode below 5 KOhm	No	Yes	Yes	Yes	

1

2

### 3 Artifact Removal and Data Preprocessing

4

5 All EEG data were imported into EEGLAB 2025.0.0 (Delorme & Makeig, 2004) and processed  
6 using 2 pipelines: a pipeline that follows the original study as closely as possible (see Vogel  
7 and Machizawa, 2004), and a recent pipeline optimized for current advances in the field of  
8 neuroscience.

9

#### 10 Direct replication preprocessing pipeline

11

12 Following the Vogel and Machizawa (2004) study, we downsampled the data to 250 Hz and  
13 applied a bandpass filter of 0.01-80 Hz (half-power cutoff, Butterworth filters) using the  
14 EEGLAB function `pop_eegfiltnew` (Widmann & Schröger, 2012). Since certain labs measured  
15 eye movements using an eye tracker or did not have access to a faraday cage, line noise (50  
16 Hz in Europe, 60 Hz in the US) was introduced as a result. To mitigate this, we used ZapLine  
17 Plus, which adaptively identifies and suppresses power-line components and their harmonics.

1 The algorithm is highly effective at removing power line artifacts while preserving non-  
2 artifactual parts of the signal (de Cheveigné, 2020; Klug & Kloosterman, 2022). This deviation  
3 from the original study was necessary to ensure accurate measurements. Afterwards, we re-  
4 referenced the data to an algebraic average of the left and right mastoids. We segmented the  
5 data from -200 to +1200 ms after the presentation of the memory array. The segments with  
6 saccadic eye movements (greater than  $1^\circ$  from the fixation cross) were excluded from further  
7 analysis using horizontal EOG channel response data from the saccade calibration task (for  
8 detailed information please refer to Exclusion Criteria). Furthermore, blinks were detected by  
9 using an amplitude threshold ( $>90$  microvolt) in the unipolar VEOG channel. In addition, a  
10 segment was marked as bad if any electrodes of interest (i.e., HEOG, P3, T5/P7, O1; P4,  
11 T6/P8, O2) had a peak-to-peak amplitude  $>200$  microvolt within one time window (i.e., bad  
12 channel criteria). Visual inspection was used to identify bad trials. Finally, a baseline correction  
13 was applied using a pre-stimulus interval of -200 to 0 ms.

14

### 15 Advanced preprocessing pipeline

16

17 In addition to following the original study's data preprocessing protocol, the data were also  
18 processed using recent advancements in neuroscience to assess the robustness of the  
19 results. First, error-prone channels were detected by the algorithms implemented in the  
20 EEGLAB plugin `clean_rawdata` ([http://sccn.ucsd.edu/wiki/Plugin\\_list\\_process](http://sccn.ucsd.edu/wiki/Plugin_list_process)) without  
21 applying automated subspace removal (ASR). An electrode was defined as error-prone when  
22 recorded data from that electrode were correlated at less than 0.85 to an estimate based on  
23 neighboring electrodes. Furthermore, an electrode was defined as error-prone if it had more  
24 line noise (i.e., 50 Hz in Europe, 60 Hz in USA) relative to its signal than all other electrodes  
25 (4 standard deviations). Finally, if an electrode had a longer flat line than 5 s, it was considered  
26 error-prone. These error-prone electrodes were automatically removed and later interpolated  
27 using a spherical spline interpolation (EEGLAB function `eeg_interp.m`). Next, data was filtered  
28 using a bandpass filter of 0.01-80 Hz (half-power cutoff, Butterworth filters). Again, we used  
29 ZapLine Plus to remove line noise. Subsequently, the data was re-referenced to an algebraic  
30 average of the left and right mastoids and segmented from -200 to +1200 ms after presentation  
31 of the memory array. To determine whether a trial was artifactual, three criteria were applied:  
32 First, we excluded trials with blinks and large saccadic eye movements. If eye-tracker  
33 recording was available, we excluded trials with blinks and eye movements larger than  $1^\circ$   
34 based on the eye-tracker. If no eye-tracker was available, we identified ocular artifacts using  
35 a tailored subject-specific amplitude threshold for the HEOG electrodes, which was obtained  
36 from the saccadic calibration task. Second, a sliding time window approach was adopted from

1 Adam et al. (2018). To identify trials containing blocking artifacts, a sliding time window of 200  
2 ms was shifted across the segments without overlap. If any time window contained 60 ms of  
3 flat line activity in any channel (i.e., range of amplitudes  $<0.1$  microvolt), the corresponding  
4 segment was marked as bad. Third, to identify trials containing large amplitude artifacts, non-  
5 overlapping sliding time windows of 12 ms were used. A segment was marked as bad if any  
6 electrode of interest (i.e., HEOG, P3, T5/P7, O1; P4, T6/P8, O2) had a peak-to-peak amplitude  
7  $>200$  microvolt within one time window. To foster scientific transparency and enable exact  
8 methodological replications and reproducibility, no visual inspection for bad trials rejection was  
9 conducted, because this decision is subjective. Finally, a baseline correction was applied  
10 using a pre-stimulus baseline interval of -200 to 0 ms.

11

## 12 CDA extraction

13

14 To remove the contribution of any VWM-unspecific, bilateral activity, the CDA was computed  
15 as a difference wave on a trial-by-trial basis by subtracting activity ipsilateral to cued items  
16 (presented left or right of screen center) from contralateral activity. The CDA amplitude was  
17 extracted from a time window of 300-900 ms after the onset of the memory array. We  
18 computed the mean CDA amplitude for each participant separately for each set size (i.e., 2, 4  
19 and 6 cued items). For the computation of the CDA, we used posterior parietal, lateral occipital  
20 and posterior temporal electrode sites (i.e., left electrode cluster: P3, T5/P7, O1; right  
21 electrode cluster: P4, T6/P8, O2). First, the difference was calculated in electrode pairs  
22 (P3/P4), (P7/P8), (O1/O2) and then averaged. The CDA was calculated on a trial-by-trial basis  
23 for all set sizes and for all trials (i.e., correct and incorrect trials). The final step was to compute  
24 the overall average CDA for each set size by averaging the CDA of the right and left cue  
25 direction of the respective set size.

26

## 27 Data Quality and Psychometric Internal Consistency

28

29 Estimates of data quality and psychometric internal consistency were reported. Data quality  
30 estimates characterize the precision of group-level ERP estimates, whereas internal  
31 consistency estimates indicate whether scores are measured reliably enough to differentiate  
32 between individuals, which is crucial for studying individual differences (Clayson, Brush, et al.,  
33 2021; Clayson & Miller, 2017; Luck et al., 2021). These metrics were reported to characterize  
34 the obtained data, but data were not excluded based on these metrics to be consistent with  
35 the procedures of the original study. Arithmetically derived estimates of the standard error of  
36 the mean were used to characterize data quality (Luck et al., 2021). These estimates

1 separately quantified the precision of CDA for each set size (2, 4, and 6 cued items) using  
2 single-trial estimates of CDA (contralateral-ipsilateral activity differences). Psychometric  
3 internal consistency estimates used generalizability theory equations to compute coefficients  
4 of dependability for difference scores (Baldwin et al., 2015; Brennan, 1992; Clayson, Baldwin,  
5 et al., 2021; Clayson, Brush, et al., 2021; Sundre, 1993). Time-window mean amplitude  
6 estimates of single-trial scores of ipsilateral and contralateral activities were used to estimate  
7 the observed group-level internal consistency of the difference scores. Dependability of  
8 contralateral-ipsilateral activity difference scores was estimated separately for each set size  
9 and data collection site using the ERP Reliability Analysis Toolbox (Clayson, Carbine, et al.,  
10 2021; Clayson & Miller, 2017). Because CDA scores were calculated as the difference  
11 between activity from different electrode sites on the same trial, residual covariances were  
12 estimated because the constituent events of the difference scores are co-occurring (Clayson,  
13 Baldwin, et al., 2021).

## 14 15 Outcome-neutral test

16  
17 To ensure that the data can test the stated hypotheses, we included quality checks and  
18 outcome-neutral tests. As an outcome-neutral test, we tested the presence of an asymmetry  
19 between contra- or ipsilateral electrode clusters time-locked to the memory array. For this, we  
20 averaged the ERPs across all set sizes and all subjects (i.e., grand averaged ERP) elicited by  
21 memory arrays that were either contra- or ipsilateral to electrode positions. A paired sample t-  
22 test for CDA between ipsilateral and contralateral sites was performed separately at each  
23 study site to verify the expected within-lab CDA experimental effect. If the t-test was significant  
24 ( $p < 0.05$ ) with more negative CDA for contralateral activity than for ipsilateral activity, then  
25 this pattern of effect justified moving forward with testing the proposed hypotheses.

## 26 27 Statistical analysis

28  
29 For all the statistical analyses, frequentist and Bayesian approaches were used. To estimate  
30 effect sizes, the statistical analyses were initially conducted for each participating lab  
31 separately. Because of the small sample size in each lab, we refrained from interpretation of  
32 the lab-specific statistics. However, the overall replication success for the project was  
33 determined based on meta-analytically pooled effect sizes, as per the defined criteria.

## 34 35 Statistical analysis for Hypothesis #1

36

1 A repeated-measures ANOVA of the CDA amplitude in the original study revealed a significant  
2 main effect for set size. Post-hoc t-tests showed significant increases in CDA amplitude for  
3 set sizes 4 and 6 compared to set size 2, with no significant differences between set sizes 4  
4 and 6. In accordance with the original study, we also conducted repeated-measures ANOVA.  
5 The significance level was set to  $p < 0.02$  uncorrected for multiple comparisons. If the ANOVA  
6 revealed a significant main effect, we further conducted post-hoc t-tests (with a significance  
7 level of  $p < 0.02$ , one-sided). We specifically tested one-sided, because we hypothesized a  
8 significant increase in CDA amplitude from arrays of 2 items per side to arrays of 4 items per  
9 side [H1.1] and 6 items per side [H1.2]. As in the original study, we adjusted the p-values with  
10 the Greenhouse-Geisser correction for nonsphericity (Jennings & Wood, 1976). If the ANOVA  
11 revealed a significant main effect, and the post-hoc t-tests show a significant increase in the  
12 CDA amplitude between arrays of 2 items per side and arrays of 4 items per side or 6 items  
13 per side, it supported hypotheses [H1.1] and [H1.2], respectively.

14  
15 We ran the corresponding analyses in a Bayesian analytical framework using a Bayesian  
16 generalized linear mixed models implemented in the brms R package (Bürkner, 2017). In the  
17 following formulas fixed effects are denoted by a "+" symbol and interaction effects by an "\*"   
18 symbol, in line with the Wilkinson notation (Wilkinson & Rogers, 1973). The predictor variable  
19 was set size (factor with 3 levels: set sizes 2, 4 and 6; reference level = set size 2). The  
20 covariates included gender (factor with 2 levels: male, female; reference level = female),  
21 handedness (factor with 2 levels: right, left; reference level = right) from Edinburgh  
22 Handedness Inventory (EHI). Due to the very small number of participants in the other and  
23 ambidextrous categories, these subjects were excluded from the Bayesian analyses. Set size  
24 and all covariates were modeled as fixed effects, while site and subject were included as  
25 random effects. Please note that the Bayesian models included only random intercepts for  
26 subject and laboratory. More complex random effects structures, such as random slopes or  
27 fully specified hierarchical models, led to frequent non convergence and unstable parameter  
28 estimates. To ensure reliable and interpretable results, we therefore adopted a parsimonious  
29 random effects specification that captured between subjects and between laboratory variability  
30 while maintaining stable model estimation. Subsequently, the credible intervals (CIs) of the  
31 posterior distributions were calculated from the newly estimated levels of significance. We  
32 opted not to calculate Bayes factors for point estimates to determine whether the effect was  
33 zero or unequal to zero. This decision was made because these Bayes factors, which rely on  
34 the Savage-Dickey ratio, heavily depend on the selection of the prior distribution for each  
35 effect. Instead, we employed a different approach: we considered a model parameter to be  
36 significant if its 98%-CI did not include zero. As suggested by Gelman (2007), the predictors  
37 and outcome variables were scaled to achieve a mean of 0 and a standard deviation of 0.5

1 (Gelman, 2007). For initial prior distributions, uninformative Cauchy priors were set to a mean  
2 of 0 and a standard deviation of 2.5.

3  
4 Importantly, the original Vogel and Machizawa (2004) study conducted ANOVA without  
5 covariates. To ensure direct compatibility with the original findings, our primary replicated  
6 analysis therefore also relied on an ANOVA without covariates. Additional analyses extending  
7 the original design conducted within the Bayesian framework included covariates to account  
8 for between subject and site-relevant variability.

## 9 10 Statistical analysis for Hypothesis #2

11  
12 In the original paper, VWM capacity was positively correlated with the CDA amplitude increase  
13 from set size 2 to 4 (i.e., when the smaller set size is below typical adult working memory  
14 capacity estimates), but not from set size 4 to 6 (i.e., when both set sizes are at or exceed  
15 capacity estimates for typical adults). To replicate this, we calculated the mean CDA amplitude  
16 increase from set size 2 to 4, and from set size 4 to 6, for each subject individually. The VWM  
17 capacity was calculated using the same formula as in the original study. This formula was  
18 introduced by Pashler (1988) and refined by Cowan (2001). It is based on the assumption that  
19 if a person can retain  $K$  items from an  $S$ -item array, then the changed item should be among  
20 the  $K$  items being held in memory on  $(K/S)$  trials, leading to correct answers on  $(K/S)$  trials  
21 where an item changed (Cowan, 2001; Pashler, 1988). The formula accounts for the false  
22 alarm rate to adjust for guessing and is expressed as  $K = S \times (H - F)$ , where  $K$  is the memory  
23 capacity,  $S$  is the set size,  $H$  is the observed hit rate in the given set size, and  $F$  is the false  
24 alarm rate in the given set size. The resulting  $K$  scores from all set sizes (i.e., 2, 4, 6) were  
25 used to compute an average  $K$  score, which we used as the behavioral measure of VWM  
26 capacity. The relationship between VWM capacity and an increase in CDA amplitude from 2  
27 to 4 items was statistically tested using Pearson's correlation. The significance level was set  
28 to  $p < 0.02$  (one-sided). If the Pearson's correlation revealed a significant positive relationship  
29 between VWM capacity and the CDA amplitude increase from 2 to 4 items, it supported the  
30 hypothesis [H2.1]. Furthermore, we conducted a Bayesian linear mixed model with a prior  
31 assuming the reported correlation coefficient from the original study ( $r = 0.78$ ). Again,  
32 significance is considered if the 98%-CI of the model parameter does not include zero.

## 33 34 Replication Success

35

1 Replication success was assessed for each hypothesis separately and was defined  
2 operationally as a statistically significant random-effects meta-analytic estimate (at  $p < 0.02$ )  
3 combining the results from the different laboratories, in the same direction as in the original  
4 study.

5  
6 Hypothesis [H1.3] and [H2.2] were analyzed using an equivalence test for meta-analyses  
7 (Lakens, 2017). The equivalence test assesses whether the difference in CDA amplitude  
8 between arrays of 4 items and 6 items is as extreme as the smallest effect size of interest  
9 (SESOI) using the two one-sided tests (TOST) procedure implemented in the R package  
10 TOSTER (Caldwell, 2022; Lakens, 2017). To perform TOST, the SESOI and its lower and  
11 upper equivalence bounds must be established. Simonsohn recommended specifying the  
12 equivalence bounds for replication studies using the “small telescopes approach” (Simonsohn  
13 et al., 2015). The idea is to consider the effect size that would give the original study 33%  
14 power. If the original study had 33% power, the probability of observing a significant effect, if  
15 there was a true effect, is too low to reliably distinguish signal from noise. Using the small  
16 telescopes approach for hypothesis [H1.3], the SESOI is  $d = 0.36$ . An alternative approach  
17 would be to calculate the smallest effect size that can be detected at a predefined power level  
18 (e.g., 90%), given the sample size and alpha level. With this approach the smallest effect size  
19 would be very similar to the small telescope approach (i.e.,  $d = 0.44$ ). Therefore, we decided  
20 to define the SESOI based on the “small telescopes approach” (i.e.,  $d = 0.36$ ) as this approach  
21 was specifically recommended for replication studies. The TOST procedure was then  
22 conducted against these bounds based on the SESOI. If the 96% confidence interval of the  
23 meta-analytic effect size falls within the equivalence bounds, the observed meta-analytic effect  
24 is statistically equivalent (Lakens, 2017). In order to test hypothesis [2.2], which postulated  
25 that there is no correlation between the subject's VWM capacity and the CDA amplitude  
26 increase from 4 to 6 items, we conducted another equivalence test. Similar to hypothesis [1.3],  
27 we used the small telescope approach to specify the SESOI (i.e.,  $r = \pm 0.29$ ).

28  
29 Finally, sequential Bayesian updating was employed by fitting a Bayesian model for each  
30 hypothesis separately to each dataset. The posterior distributions obtained from each analysis  
31 were used as priors for the next analysis, allowing evidence to be accumulated across the  
32 datasets from different labs. This approach was expected to produce greater statistical power  
33 than independent analyses and yield more robust outcome parameters.

34

35 **Sensitivity analyses**

36

1 Recently, there have been concerns regarding the validity of the K score as a measure of  
2 VWM capacity. Specifically, some researchers in the field have noted that K operates under  
3 the assumption of all-or-none memories and does not account for individual decision biases,  
4 which can lead to an overestimation of capacity depending on the observer's strategy (Brady  
5 et al., 2023; Williams et al., 2022).

6  
7 In light of these concerns, we conducted an additional analysis using the  $d'$  ( $d$  prime) metric.  
8  $d'$  is a commonly used measure in signal detection theory and provides a unitless, normalized  
9 measure of sensitivity that is independent of response bias.  $d'$  is defined as  $d' = Z(\text{hit rate}) - Z(\text{false alarm rate})$ . A hit was defined as reporting a color change when there was one,  
10 and a false alarm as reporting such a change when no change occurred. The resulting  $d'$   
11 scores from all set sizes (i.e., 2, 4, 6) were used to compute an average  $d'$  score. Our  
12 reasoning is that replicating the results using both K score (as the primary analysis) and  $d'$  (as  
13 an additional analysis) would provide stronger evidence for the observed effect, as it would  
14 demonstrate that the results are not solely dependent on the characteristics of the K measure.  
15

## 16 17 Deviations from Preregistration

18 A small number of deviations from the preregistered Stage-1 protocol were necessary to  
19 ensure acceptable data quality across laboratories.

20  
21 (1) We preregistered that data would be collected across 10 laboratories. However, after Stage  
22 1 acceptance, the laboratory in Finland withdrew from the project due to regulatory reasons.  
23 After the laboratory in Finland withdrew from the project, several participating labs agreed to  
24 recruit additional participants to maintain the planned statistical power. Subsequently, we were  
25 also able to include an additional site (Prof. Sauseng's lab from Zurich; UZH-NCN), which  
26 further increased the total sample size. As a result, the total number of participants exceeded  
27 the preregistered target and was unevenly distributed across laboratories (Table 3).

28  
29 (2) The preregistered blink threshold of 50  $\mu\text{V}$  in the VEOG channel resulted in an excessive  
30 number of rejected trials in several laboratories. In the seminal CDA study by Vogel and  
31 Machizawa (2004), the authors state that "ERPs were recorded using our standard recording  
32 and analysis procedures, including rejection of trials contaminated by blinks or large eye  
33 movements," citing earlier work from 1998 (Vogel et al., 1998). However, neither the 2004  
34 paper nor the referenced 1998 work specifies an explicit VEOG amplitude threshold for blink  
35 detection. Similarly, other early CDA studies do not provide a clearly defined numeric criterion  
36 for blink rejection based on VEOG amplitude (Drew & Vogel, 2008; Fukuda & Vogel, 2011;

1 McCollough et al., 2007; Vogel & Machizawa, 2004). In later work from the Vogel/Awh lab,  
2 blink detection was implemented using algorithmic, sliding-window step-function approaches  
3 applied to the VEOG signal. For example, Adam (2018) detected blinks using a sliding window  
4 on VEOG (window size = 200 ms, step size = 10 ms, threshold = 50  $\mu$ V). Subsequent studies  
5 (Hakim et al., 2021) further refined this approach by using smaller thresholds (e.g., 30  $\mu$ V;  
6 window size = 80 ms, step size = 10 ms). Importantly, these algorithmic thresholds are not  
7 directly comparable to our single-value VEOG amplitude cutoff, as sliding-window methods  
8 detect rapid step-like deflections over short temporal windows rather than absolute VEOG  
9 magnitude at a single time point. In the present study, the preregistered blink criterion of 50  
10  $\mu$ V was applied as a single-value VEOG threshold, which resulted in an excessive number of  
11 rejected epochs in several laboratories. Because single-value approaches require higher  
12 cutoffs to achieve a comparable level of stringency, we increased the blink threshold to 90  $\mu$ V  
13 to retain an adequate number of usable epochs across laboratories.

14

15 (3) We used ZaplinePlus instead of the preregistered ZapLine algorithm. Several laboratories  
16 exhibited substantial line-noise contamination, and ZaplinePlus provided more robust  
17 identification and suppression of line-noise harmonics while preserving neural signals.

18

19 (4) The preregistered criterion defined blocking as amplitude ranges  $<1$   $\mu$ V for 30 ms. Applying  
20 this threshold would have resulted in the rejection of an unacceptably large proportion of trials  
21 ( $>30\%$  on average). The threshold was therefore adjusted to  $<0.1$   $\mu$ V for 60 ms, which  
22 preserved data quality while avoiding disproportionate data loss and is in line with the logic of  
23 flatline detection implemented in automated preprocessing pipelines such as `clean_rawdata`  
24 (Kothe & Makeig, 2013).

25

26 (5) In a subset of laboratories (USF, OSU, UZH-NCN), the eye calibration task did not yield  
27 accurate HEOG amplitudes, resulting in estimates that were either unrealistically high or too  
28 low and therefore would have led to retaining too many trials or removing an excessive number  
29 of trials. For these laboratories, we applied a fixed HEOG threshold of 30  $\mu$ V, which closely  
30 matched the heuristic value used in the seminal Vogel and Machizawa (2004) and the average  
31  $1^\circ$  amplitude obtained from the laboratories with valid eye calibration data.

32

33 (6) The preregistered plan stated that equivalence testing would be evaluated using 90%  
34 confidence intervals, consistent with the standard TOST framework ( $\alpha = 0.05$ ). However, all  
35 EEGManyLabs replication studies adopt a project-wide significance threshold of  $\alpha = 0.02$  to  
36 ensure a uniform inferential standard across analyses and tasks. The threshold of  $\alpha = 0.02$   
37 was applied consistently throughout the present project for all frequentist statistical tests. To

1 maintain internal consistency with the EEGManyLabs inferential framework, the equivalence  
2 tests were therefore conducted using 96% confidence intervals (i.e.,  $1 - 2\alpha$ ), rather than the  
3 preregistered 90%.

4

5 (7) An exploratory preprocessing pipeline incorporating ICA followed by ICLabel-based  
6 component classification was implemented to assess robustness of the results. The ICA  
7 pipeline was not preregistered and is therefore reported only in the Supplement.

8

Journal Pre-proof

### 1 3. Results

2

3 To provide a concise overview of the results across hypotheses, preprocessing pipelines, and  
 4 statistical frameworks, we summarize the replication outcomes in Table 3. Across the direct,  
 5 advanced, and ICA pipelines, outcome neutral effects and set-size-related hypotheses (H1.1-  
 6 H1.3) were consistently replicated using both frequentist and Bayesian approaches. In  
 7 contrast, the hypothesis targeting the association between the CDA increase from set size 2  
 8 to 4 and individual VWM capacity (H2.1) showed less consistent support across analyses, with  
 9 replication depending on the specific pipeline and modelling approach and effects consistently  
 10 remaining close to the threshold of statistical significance. Finally, equivalence between the  
 11 CDA increase from set size 4 to 6 and individual VWM capacity (H2.2) was again consistently  
 12 supported across all analyses.

13

14 Detailed results for the advanced and ICA pipelines are reported in the Supplementary  
 15 Materials, as their outcomes closely mirrored those of the direct pipeline and are included  
 16 there to reduce redundancy and streamline the presentation of the results.

17

18 Table 3. Summary of the results across analyses and preprocessing pipelines

	Direct			Advanced			ICA		
	Frequen- tist	Bayesian LMM	Bayesian Seq.	Frequen- tist	Bayesian LMM	Bayesian Seq.	Frequen- tist	Bayesian LMM	Bayesian Seq.
Outcome N.	✓	✓	✓	✓	✓	✓	✓	✓	✓
H1.1.	✓	✓	✓	✓	✓	✓	✓	✓	✓
H1.2.	✓	✓	✓	✓	✓	✓	✓	✓	✓
H1.3.	✓	✓	✓	✓	✓	✓	✓	✓	✓
H2.1.	X	X	X	X	X	X	X	✓	✓
H2.2.	✓	✓	✓	✓	✓	✓	✓	✓	✓

19 Note. ✓ = indicates successful replication. X = indicates failure to replicate. LMM = Linear Mixed Model.  
 20 Seq = Sequential Updating.

21

#### 22 3. 1. Sample Characteristics and Artifact-Related Trial Rejection

23

24 Table 4 provides an overview of the demographic characteristics of all participants across the  
 25 10 laboratories. In total, 304 participants were recruited. After trial rejection and subject

1 exclusion, the final sample comprised 217 participants in the direct pipeline, 231 participants  
 2 in the advanced pipeline, and 266 participants in the ICA pipeline. For comparison, the original  
 3 study (Vogel & Machizawa, 2004) reported only limited demographic information, including  
 4 the total sample size (N = 36) and an age range of 21-33 years. No further details (e.g., sex,  
 5 handedness) were provided.

6

7 Table 4. Demographic characteristics of the full sample (N = 304).

	N	Age		Gender			Handedness		
		M	SD	Female	Male	Other	R	L	A
DART	28	20.29	3.41	9	18	1	23	3	0
USF	50	19.64	1.90	41	9	0	34	11	3
JGU	25	23.04	3.02	17	8	0	22	2	1
WWU	30	22.80	2.27	25	5	0	28	2	0
NDSU	28	24.93	5.80	15	13	0	23	5	0
OSU	24	24.13	3.39	18	6	0	NA	NA	NA
UI	30	24.50	4.90	21	7	0	25	4	0
TUOS	25	20.04	2.95	20	5	0	23	2	0
UZH-MPR	40	23.35	2.82	25	15	0	40	0	0
UZH-NCN	24	22.92	3.97	16	8	0	24	0	0

8 Note. M = Mean. SD = Standard deviation. R = Right. L = Left. A = Ambidextrous.

9

10 Table 5 summarizes the number of rejected trials across laboratories for each exclusion  
 11 criterion and preprocessing pipeline. As preregistered, trials were rejected due to amplitude  
 12 thresholds, blocking artefacts, VEOG artefacts (i.e., blinks), and HEOG artefacts (i.e.,  
 13 saccades). In addition, 4 laboratories collected eye-tracking data. For these sites, blink and  
 14 saccade rejections in the advanced and ICA pipelines were based on ET-detected events  
 15 rather than VEOG/HEOG thresholds. In contrast, the original study (Vogel & Machizawa,  
 16 2004) only reports that trials containing artefacts were rejected, without providing quantitative  
 17 information on the number or proportion of rejected trials.

18

19 Table 5. Number of rejected trials by rejection criterion, laboratory, and preprocessing pipeline.

Pipeline	Amplitude		Blocking		VEOG		HEOG		ET Blink		ET Saccade		Final N
	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	
Dartmouth College													
Direct	48.18	109.06	0	0	76.79	122.39	81.89	126.94	-	-	-	-	23
Advanced	10.14	27.11	0	0	32.21	35.01	38.64	83.92	-	-	-	-	27
ICA	10.18	27.17	0	0	31.86	35.47	36.75	81.7	-	-	-	-	27
University of South Florida													
Direct	99.14	130.42	0	0	165.24	151.02	97.44	129.18	-	-	-	-	26
Advanced	67.1	99.34	0	0	134.58	143.48	71.4	91.71	-	-	-	-	33
ICA	60.22	96.66	0	0	83.64	150.88	40.56	75.12	-	-	-	-	40
Johannes Gutenberg University Mainz													
Direct	3.56	6.93	0	0	129.12	152.28	174.96	123.82	-	-	-	-	11
Advanced	3.52	6.95	0	0	128.68	152.34	174.84	123.76	-	-	-	-	11
ICA	1.36	2.97	0	0	3.72	7.23	4.52	11.72	-	-	-	-	25
University of Münster													
Direct	40.27	98.91	0.1	0.4	34.27	36.03	36.27	46.43	-	-	-	-	26
Advanced	38.47	134.37	0.43	2.03	36.67	43.52	35.27	93.58	21.97	22.33	19.53	20.08	28
ICA	32.83	132.98	0.47	2.37	28.2	130.24	24.2	131.42	21.97	22.33	19.53	20.08	28
North Dakota State University													
Direct	5.29	11.44	0	0	78.93	97.46	63.68	78.58	-	-	-	-	22
Advanced	2.57	6.81	0	0	65.89	98.55	63.82	78.82	49.18	74.41	52.43	85.41	25
ICA	2.25	6.54	0	0	6.86	12.75	6.32	23.76	49.18	74.41	52.43	85.41	25
The Ohio State University													
Direct	15.71	32.92	0	0	76.04	89.99	26.21	59.38	-	-	-	-	20
Advanced	8.88	14.74	0	0	35.21	61.92	13.96	17.78	82.42	94.46	100.25	109.43	19
ICA	8.21	14.69	0	0	6.58	15.08	11.33	18.89	82.42	94.46	100.25	109.43	19
University of Iceland													
Direct	11.8	32.44	0	0	79.83	121	79.6	82.64	-	-	-	-	24
Advanced	11.37	31.91	0	0	6.57	15.96	144.67	195.85	-	-	-	-	22
ICA	9.77	31.05	0	0	4.47	14.74	86.87	165.68	-	-	-	-	24
University of Sheffield													

Direct	11.92	24.8	0	0	135.6	127.05	106.96	93.48	—	—	—	—	15
Advanced	11.52	23.68	0	0	135.68	127.77	106.76	93.47	—	—	—	—	16
ICA	10.68	23.74	0	0	32.68	61.54	8.96	16.66	—	—	—	—	24
University of Zurich (UZH-MPR)													
Direct	41.1	88.56	0	0	103.98	131.67	86.98	114.74	—	—	—	—	30
Advanced	2.75	7.53	0	0	66.65	105.35	86.83	114.73	46.88	106.69	148.45	163.04	30
ICA	2.3	6.35	0	0	9.1	19.83	44.68	152.02	46.88	106.69	148.45	163.04	30
University of Zurich (UZH-NCN)													
Direct	11.33	25.82	0	0	93.17	137.66	119.71	115.23	—	—	—	—	20
Advanced	11.21	25.74	0	0	87.42	137.2	111.71	119.02	—	—	—	—	20
ICA	1.63	3.31	0	0	3.58	8.64	22.58	25.69	—	—	—	—	24

1 Note. M = Mean. SD = Standard deviation. Eye-tracking (ET) data and ET-based rejection criteria were  
 2 available only for the WWU, NDSU, OSU, and UZH-MPR laboratories. An em dash (—) indicates that  
 3 ET data were either not recorded at that site or were not included in the analysis pipeline.

4

### 5 3. 2. Behavioral Results

6

7 Performance in the color change detection task, expressed as accuracy, was  $87.7\% \pm 14.5\%$   
 8 for set size 2,  $75.8\% \pm 13.0\%$  for set size 4, and  $65.2\% \pm 10.5\%$  for set size 6 (Table 6). The  
 9 average accuracy across all set sizes was  $76.2\% \pm 12.0\%$ . Performance in the color change  
 10 detection task, expressed as K scores, was  $1.61 \pm 0.52$  for set size 2,  $2.26 \pm 0.93$  for set size  
 11 4, and  $2.11 \pm 1.10$  for set size 6. The average K score across all set sizes was  $1.99 \pm 0.80$ .  
 12 Performance expressed in  $d'$  (d-prime) was  $3.09 \pm 1.26$  for set size 2,  $1.82 \pm 0.87$  for set size  
 13 4, and  $1.14 \pm 0.62$  for set size 6, with an overall mean  $d'$  of  $2.01 \pm 0.86$ . The average K score  
 14 and average  $d'$  were strongly correlated ( $r = 0.93$ ,  $p = 2.93e-95$ ), indicating high convergence  
 15 between the two behavioral measures. Additionally, all further analyses based on K-score and  
 16  $d'$  yielded highly similar results. The original study (Vogel & Machizawa, 2004) did not report  
 17 detailed behavioral summary statistics (e.g., accuracy, K-scores, or  $d'$ ), although an  
 18 approximate average K-score of 2.8-2.9, averaged across set sizes 2, 4, and 6, can be visually  
 19 inferred from the published figure.

20

21 Table 6. Performance in the color change detection task.

Accuracy (%)				K-score				D-prime			
2	4	6	Average	2	4	6	Average	2	4	6	Average

	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD	M	SD
DART	91	6	80	8	69	8	80	11	1.71	0.18	2.53	0.56	2.48	0.87	2.24	0.71	3.24	0.88	2.00	0.58	1.29	0.46	2.17	1.04		
USF	75	23	65	18	56	12	65	20	1.27	0.89	1.67	1.25	1.41	1.13	1.45	1.11	2.19	1.69	1.26	1.02	0.73	0.65	1.39	1.33		
JGU	90	10	77	9	67	8	78	13	1.67	0.37	2.36	0.78	2.23	1.03	2.09	0.82	3.25	1.01	1.89	0.75	1.2	0.55	2.11	1.16		
WWU	94	10	85	9	73	8	84	13	1.81	0.39	2.93	0.72	2.97	1.02	2.57	0.92	3.99	1.07	2.48	0.73	1.6	0.6	2.69	§.28		
NDSU	84	12	68	10	59	8	70	14	1.5	0.46	1.72	0.79	1.41	0.84	1.54	0.72	2.76	1.23	1.32	0.62	0.74	0.41	1.61	1.19		
OSU	93	5	82	6	70	7	82	11	1.76	0.15	2.68	0.46	2.46	0.83	2.30	0.68	3.41	0.78	2.19	0.66	1.31	0.45	2.31	1.07		
UI	86	16	72	11	61	7	73	16	1.49	0.65	1.83	0.84	1.5	0.77	1.61	0.77	2.75	1.36	1.43	0.69	0.83	0.43	1.67	1.21		
TUOS	88	8	74	11	64	10	75	14	1.59	0.25	2.06	0.84	2.03	1.12	1.89	0.84	2.78	0.7	1.68	0.93	1.18	0.76	1.88	1.04		
UZH-MPR	92	8	80	9	69	8	80	13	1.75	0.31	2.55	0.67	2.43	0.9	2.24	0.76	3.5	0.91	2.14	0.68	1.35	0.54	2.33	1.15		
UZH-NCN	92	7	81	9	70	9	81	12	1.77	0.2	2.65	0.66	2.66	0.93	2.36	0.78	3.58	0.96	2.2	0.72	1.44	0.54	2.41	1.16		

1 Note. M = Mean. SD = Standard deviation.

2

### 3 3. 3. Data Quality and Psychometric Internal Consistency

4

5 Estimates of data quality (i.e., standardized measurement error; SME) and psychometric  
6 internal consistency (i.e., dependability coefficients) were computed for each laboratory and  
7 each set size for all 3 preprocessing pipelines (Clayson, Brush, et al., 2021). Dependability  
8 coefficients for the CDA, derived from generalizability theory and conceptually analogous to  
9 coefficient alpha in classical test theory (Clayson, Carbine, et al., 2021; Shavelson & Webb,  
10 2012) were generally high across sites, with most laboratories achieving values above 0.75  
11 (Table 7). For set size 2, dependability ranged from 0.39 to 0.90. A similar pattern was  
12 observed for set size 4, with coefficients spanning 0.36 to 0.93. For set size 6, dependability  
13 values ranged from 0.35 to 0.92. Overall, CDA dependability coefficients in the present study  
14 were generally within or above the range considered acceptable for preliminary research ( $\geq$   
15 0.70) and, for many laboratories, approached or exceeded the more stringent threshold  
16 recommended for studies of group differences ( $\geq$  0.80; Clayson & Miller, 2017; Nunnally &  
17 Bernstein, 1994).

18

19 Data quality, indexed by the SME of single-trial CDA amplitudes, showed comparable patterns  
20 across sites and set sizes (Table 7). SME values for set size 2 ranged from 0.297 to 1.439.  
21 For set size 4, SME values ranged from 0.299 to 1.382, and for set size 6, values ranged from  
22 0.297 to 1.426. Across the majority of laboratories, SME values clustered between  
23 approximately 0.35-0.50, indicating precise trial-level CDA estimates. USF again showed  
24 substantially larger SME values across all set sizes, consistent with its lower dependability  
25 coefficients. Importantly, as emphasized by Luck and colleagues (2021), SME is intended as

1 a continuous index of data quality rather than a metric with predefined acceptability thresholds,  
 2 and it is therefore difficult to classify SME values as intrinsically “good” or “bad”. Instead, lower  
 3 SME values indicate higher precision and reliability. Together, these results demonstrate that  
 4 CDA amplitudes were estimated with good precision and strong internal consistency across  
 5 most participating laboratories, providing confidence in the reliability of the measurements  
 6 underlying the primary analyses.

7

8 Table 7. Data Quality Metrics Across Laboratories.

Pipeline	Dependability (Estimate and 98%-CI)			SME		
	Set size 2	Set size 4	Set size 6	Set size 2	Set size 4	Set size 6
Dartmouth College						
Direct	0.84 (0.71, 0.92)	0.88 (0.77, 0.94)	0.87 (0.76, 0.94)	0.36	0.35	0.38
Advanced	0.83 (0.69, 0.92)	0.88 (0.78, 0.94)	0.88 (0.78, 0.94)	0.36	0.34	0.36
ICA	0.83 (0.69, 0.92)	0.88 (0.78, 0.94)	0.87 (0.77, 0.94)	0.36	0.34	0.36
University of South Florida						
Direct	0.39 (0.18, 0.60)	0.36 (0.15, 0.58)	0.35 (0.15, 0.57)	1.44	1.38	1.43
Advanced	0.51 (0.31, 0.70)	0.50 (0.28, 0.69)	0.51 (0.29, 0.69)	1.22	1.16	1.24
ICA	0.42 (0.21, 0.62)	0.51 (0.29, 0.70)	0.43 (0.22, 0.64)	1.30	1.20	1.16
Johannes Gutenberg University Mainz						
Direct	0.85 (0.70, 0.93)	0.87 (0.75, 0.94)	0.89 (0.78, 0.95)	0.38	0.38	0.35
Advanced	0.84 (0.69, 0.93)	0.87 (0.75, 0.94)	0.88 (0.77, 0.95)	0.38	0.38	0.35
ICA	0.85 (0.71, 0.93)	0.88 (0.76, 0.94)	0.90 (0.80, 0.95)	0.39	0.27	0.28
University of Münster						
Direct	0.78 (0.63, 0.89)	0.71 (0.50, 0.85)	0.80 (0.65, 0.90)	0.45	0.47	0.46
Advanced	0.79 (0.63, 0.89)	0.74 (0.56, 0.87)	0.81 (0.67, 0.91)	0.40	0.40	0.41
ICA	0.77 (0.61, 0.89)	0.78 (0.62, 0.89)	0.86 (0.74, 0.93)	0.41	0.40	0.41
North Dakota State University						
Direct	0.77 (0.60, 0.89)	0.77 (0.59, 0.89)	0.79 (0.63, 0.90)	0.40	0.40	0.39
Advanced	0.80 (0.64, 0.90)	0.75 (0.57, 0.88)	0.82 (0.67, 0.91)	0.38	0.39	0.38
ICA	0.77 (0.59, 0.88)	0.74 (0.53, 0.87)	0.80 (0.63, 0.90)	0.36	0.38	0.37
The Ohio State University						

Direct	0.76 (0.56, 0.89)	0.84 (0.69, 0.92)	0.86 (0.75, 0.94)	0.39	0.39	0.41
Advanced	0.75 (0.54, 0.88)	0.84 (0.71, 0.92)	0.86 (0.75, 0.94)	0.42	0.43	0.44
ICA	0.71 (0.48, 0.87)	0.84 (0.70, 0.92)	0.83 (0.69, 0.92)	0.44	0.44	0.46
University of Iceland						
Direct	0.81 (0.67, 0.90)	0.84 (0.71, 0.92)	0.77 (0.62, 0.88)	0.42	0.42	0.43
Advanced	0.82 (0.68, 0.91)	0.84 (0.71, 0.92)	0.78 (0.60, 0.89)	0.47	0.49	0.57
ICA	0.81 (0.67, 0.90)	0.83 (0.70, 0.92)	0.79 (0.64, 0.89)	0.44	0.44	0.49
University of Sheffield						
Direct	0.77 (0.57, 0.89)	0.80 (0.64, 0.91)	0.82 (0.66, 0.91)	0.50	0.49	0.48
Advanced	0.77 (0.57, 0.89)	0.80 (0.64, 0.90)	0.81 (0.64, 0.91)	0.49	0.49	0.48
ICA	0.85 (0.72, 0.93)	0.82 (0.66, 0.92)	0.78 (0.60, 0.90)	0.39	0.38	0.40
University of Zurich (UZH-MPR)						
Direct	0.81 (0.68, 0.89)	0.91 (0.85, 0.95)	0.87 (0.78, 0.92)	0.30	0.30	0.30
Advanced	0.85 (0.74, 0.92)	0.91 (0.85, 0.95)	0.86 (0.76, 0.92)	0.33	0.33	0.34
ICA	0.86 (0.76, 0.92)	0.91 (0.85, 0.95)	0.85 (0.76, 0.92)	0.33	0.34	0.35
University of Zurich (UZH-NCN)						
Direct	0.90 (0.81, 0.95)	0.93 (0.86, 0.96)	0.92 (0.84, 0.96)	0.50	0.36	0.70
Advanced	0.90 (0.80, 0.95)	0.93 (0.86, 0.97)	0.92 (0.84, 0.96)	0.50	0.36	0.70
ICA	0.88 (0.79, 0.95)	0.92 (0.85, 0.96)	0.91 (0.84, 0.96)	0.31	0.31	0.31

Note. CI = Confidence Interval. SME = Standardized Measurement Error.

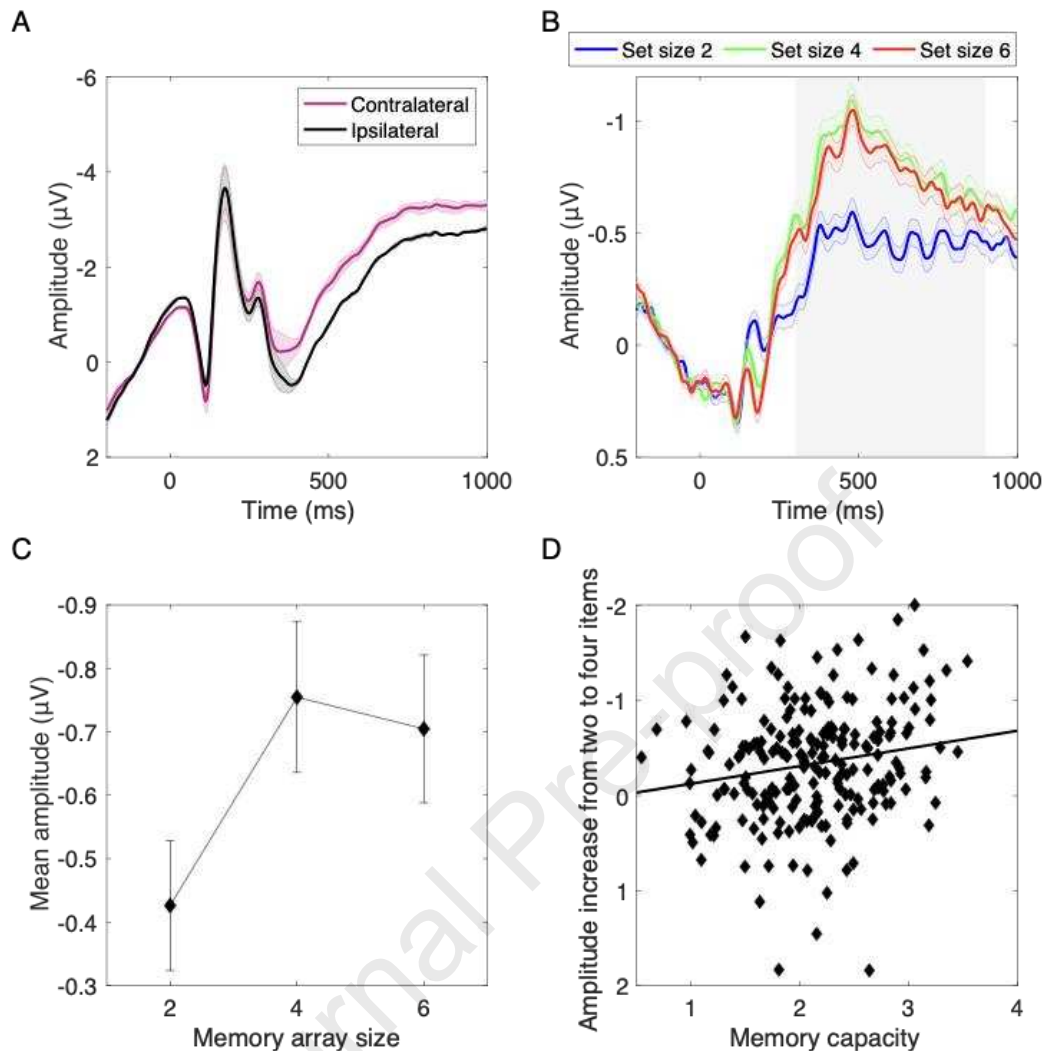
2

### 3 3. 4. Results of the Direct Replication

4

5 Here, we report the results of the direct replication pipeline, which follows the procedures of  
6 the original study as closely as possible. Figure 3 displays the grand-average contralateral  
7 and ipsilateral ERPs, the grand-average CDA with set size effects, an error-bar plot with 98%-  
8 CI illustrating the set size effects, as well as the relationship between memory capacity defined  
9 as K-score and the amplitude increase from 2 to 4 items, with all subjects from all laboratories  
10 collapsed into a single combined dataset.

11



1  
 2 Figure 3. Direct Pipeline. Grand-average CDA effects across all participants and labs, averaged across  
 3 posterior electrode clusters (left: P3, P7, O1; right: P4, P8, O2). (A) Outcome Neutral Test: Contralateral  
 4 vs. ipsilateral ERP. Subplots (B) and (C) illustrate the set size effect specified in hypothesis #1, with (B)  
 5 showing ERP responses by set size and (C) displaying CDA amplitudes for set sizes 2, 4, and 6 with  
 6 98% confidence intervals. The grey shaded region in (B) indicates the time window used for statistical  
 7 analysis. (D) Association between the CDA increase from set size 2 to 4 and VWM capacity (K-score).  
 8 The line represents the best linear fit. For visualization purposes only, participants were pooled across  
 9 laboratories into a single combined sample.

10

### 11 3. 4. 1. Outcome Neutral Test

12

13 We first conducted the outcome-neutral test to assess whether the expected contralateral-  
 14 ipsilateral asymmetry was present across labs, using a frequentist approach, a Bayesian  
 15 generalized linear mixed-model and a Bayesian sequential updating procedure.

16

### 3. 4. 1. 1. Frequentist Approach

2

3 The frequentist approach revealed a robust contralateral-ipsilateral asymmetry across  
 4 laboratories. Paired-sample t-tests at each lab consistently yielded significant contralateral-  
 5 ipsilateral differences (all  $p < 0.05$ ), with contralateral activity exhibiting more negative  
 6 amplitudes than ipsilateral activity (Table 8).

7

8 Table 8. Direct Pipeline. Paired-sample t-tests showing contralateral-ipsilateral CDA asymmetry at each  
 9 laboratory.

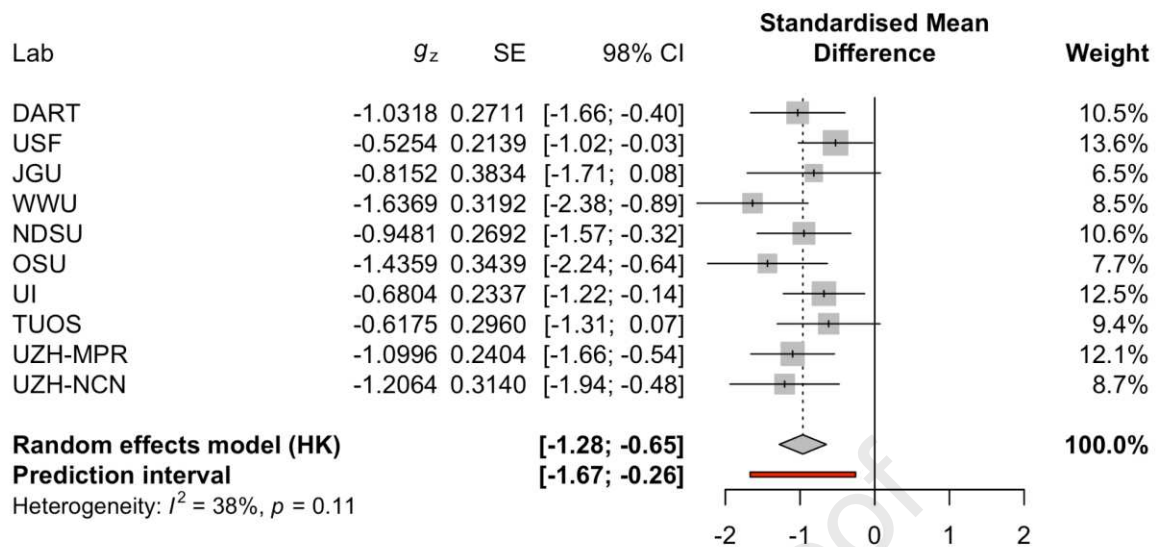
Lab	M <sub>Contra</sub>	SD <sub>Contra</sub>	M <sub>Ipsi</sub>	SD <sub>Ipsi</sub>	t-value	df	p-value	Hedges' g <sub>z</sub>
DART	-2.22	1.56	-1.48	1.69	-5.13	22	3.89e-05***	-1.03
USF	-0.33	1.02	-0.07	0.98	-2.76	25	0.011*	-0.53
JGU	-2.91	1.35	-2.42	1.22	-2.93	10	0.015*	-0.82
WWU	-2.66	2.00	-1.76	1.82	-8.61	25	6.02e-09***	-1.64
NDSU	-1.68	1.58	-1.16	1.44	-4.61	21	1.50e-04***	-0.95
OSU	-1.88	1.56	-1.00	1.16	-6.69	19	2.14e-06***	-1.44
UI	-2.94	1.54	-2.50	1.29	-3.45	23	0.002**	-0.68
TUOS	-2.35	1.41	-1.91	1.19	-2.53	14	0.024*	-0.62
UZH-MPR	-1.75	1.59	-1.09	1.57	-6.18	29	9.61e-07***	-1.10
UZH-NCN	-1.75	1.31	-0.84	1.09	-5.62	19	2.03e-05***	-1.21

10 Note. M = Mean. SD = Standard deviation. Df = Degrees of freedom.

11 \* $p < 0.05$ . \*\* $p < 0.01$ . \*\*\* $p < 0.001$

12

13 A random-effects meta-analysis confirmed a significant overall asymmetry (average  $g_z = -0.96$ ,  
 14 98%-CI = [-1.28, -0.65],  $t = -8.62$ ,  $p = 1.21e-5$ ), replicating the results from the original study  
 15 (Figure 4). Taken together, these results verify that the dataset provides a reliable  
 16 contralateral-ipsilateral asymmetry effect, thereby fulfilling the preregistered criterion for  
 17 proceeding with the main analyses.



1

2 Figure 4. Direct Pipeline. Forest plot of the meta-analysis for outcome-neutral test: contralateral-  
 3 ipsilateral asymmetry across labs. For each laboratory, the plot shows the effect size estimate Hedges'  
 4  $g$  with its standard error and corresponding 98% confidence interval. Squares represent lab-specific  
 5 effect size estimates, with square size proportional to the inverse-variance weight, indicating the relative  
 6 contribution of each laboratory to the pooled estimate (larger squares reflect greater weight due to  
 7 higher precision). Horizontal lines denote 98% confidence intervals. The random-effects meta-analytic  
 8 estimate (Hartung-Knapp adjustment) is shown as a diamond with its 98% confidence interval, together  
 9 with the 98% prediction interval. Between-laboratory heterogeneity is quantified using  $I^2$ , with the  
 10 associated  $p$ -value reported.

11

### 12 3. 4. 1. 1. Bayesian Approach

13

14 To complement the frequentist analyses, we estimated a Bayesian generalized linear mixed  
 15 model predicting lateralized ERP amplitude from laterality (factor with 2 levels: contralateral,  
 16 ipsilateral; reference level = contralateral) while including laboratory and subject as random  
 17 intercepts. As preregistered (see Stage 1 Registered Report), all Bayesian models additionally  
 18 controlled for gender and handedness.

19

$$20 \text{ ERP Amplitude} \sim \text{Laterality} * \text{Gender} * \text{Handedness} + (1|\text{Subject}) + (1|\text{Lab})$$

21

22 CDA amplitudes were credibly more positive at ipsilateral than contralateral electrodes  
 23 (Estimate = 0.20, 98%-CI = [0.15, 0.24]), with the 98% credible interval for the laterality effect  
 24 excluding zero, indicating robust evidence for the contralateral-ipsilateral asymmetry effect  
 25 across laboratories. None of the covariates showed credible associations with CDA amplitude,  
 26 as the 98% credible intervals for gender and handedness all included zero.

1  
 2 Finally, we implemented the Bayesian sequential updating procedure to evaluate the stability  
 3 of the contralateral-ipsilateral asymmetry as evidence accumulated across laboratories. For  
 4 the first site, we fit a Bayesian regression model predicting CDA amplitude from laterality using  
 5 weakly informative (uniform) priors. The posterior from this model was then used as the prior  
 6 for the next lab, and this process was repeated iteratively across all ten sites, allowing the  
 7 posterior distribution to be continuously updated as new data were incorporated. The final  
 8 posterior distribution after integrating evidence across all labs provided clear support for the  
 9 predicted contralateral-ipsilateral asymmetry (Estimate = 0.21, 98%-CI = [0.18, 0.24]), with  
 10 CDA amplitudes more positive at ipsilateral than contralateral electrodes.

11

### 12 3. 4. 2. Hypothesis #1: Set Size Effects

13

14 Having established that all laboratories reproduced the expected contralateral-ipsilateral  
 15 asymmetry, we next turned to the analyses examining set-size-dependent changes in CDA  
 16 amplitude. In line with the preregistration, we assessed these effects using both a frequentist  
 17 repeated-measures ANOVA with post-hoc t-tests, a Bayesian generalized linear mixed-model  
 18 approach and a Bayesian sequential updating procedure.

19

#### 20 3. 4. 2. 1. Frequentist Approach

21

22 We first examined the set-size effect between arrays of 2, 4, and 6 items using the frequentist  
 23 approach. A repeated-measures ANOVA was conducted separately for each laboratory (Table  
 24 9). Out of the 10 participating sites, 6 labs (DART, JGU, WWU, NDSU, and both UZH labs)  
 25 showed a statistically significant main effect of set size (all  $p < 0.02$ ), whereas 4 labs (USF,  
 26 OSU, UI and TUOS) did not reach statistical significance ( $p > 0.02$ ).

27

28 Table 9. Direct Pipeline. Results of the repeated-measures ANOVA assessing the set size effect across  
 29 labs.

Lab	df	F-value	p-value
DART	(2, 44)	5.69	0.006**
USF	(2, 50)	0.80	0.455
JGU	(2, 20)	6.72	0.006**
WWU	(2, 50)	14.37	1.51e-5***
NDSU	(2, 42)	4.34	0.019*
OSU	(2, 38)	4.08	0.025
UI	(2, 46)	1.10	0.332
TUOS	(2, 28)	2.13	0.138
UZH-MPR	(2, 58)	31.54	5.28e-10***

UZH-NCN (2, 38) 13.09 4.74e-5\*\*\*

1 Note. Df = Degrees of freedom.

2 \* $p < 0.02$ . \*\* $p < 0.01$ . \*\*\* $p < 0.001$

3

4 Hypothesis #1.1. CDA amplitude increase from set size 2 to 4

5 Next, we tested Hypothesis H1.1, which predicted an increase in CDA amplitude from set size  
6 2 to 4. To this end, we computed post-hoc paired  $t$ -tests for all sites (Table 10). 6 labs (DART,  
7 JGU, WWU, OSU, and both UZH labs) showed statistically significant differences between set  
8 sizes 2 and 4, whereas in 4 labs the difference did not reach statistical significance (USF,  
9 NDSU, UI and TUOS).

10

11 Table 10. Direct Pipeline. Post-hoc  $t$ -tests for the set size 2 vs. 4 comparisons across laboratories.

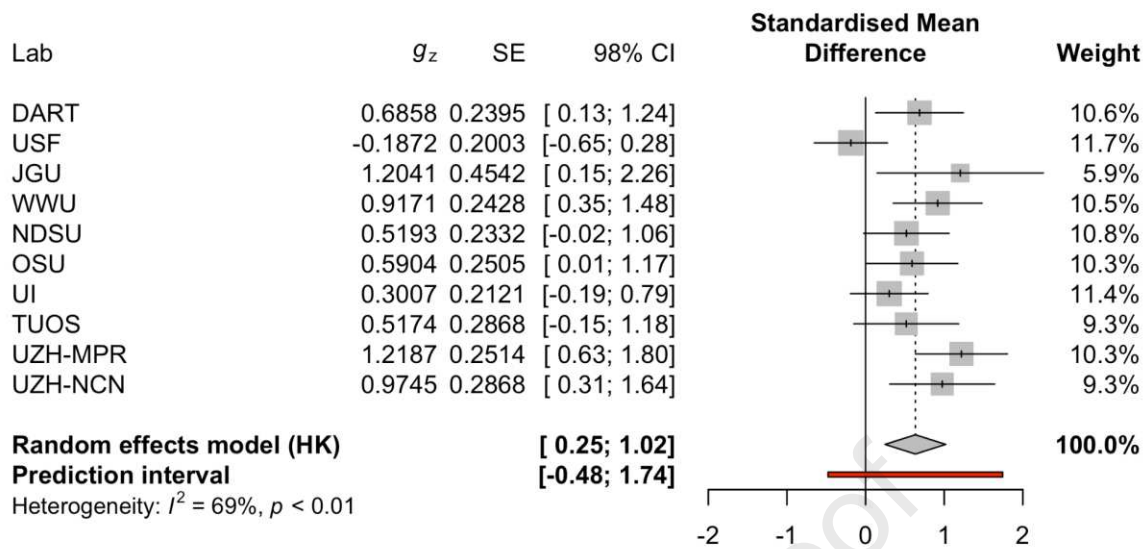
Lab	$M_{S_{22}}$	$SD_{S_{22}}$	$M_{S_{24}}$	$SD_{S_{24}}$	t-value	df	p-value	Hedges' $g_z$
DART	-0.61	0.62	-0.87	0.77	3.41	22	0.003**	0.69
USF	-0.31	0.75	-0.14	0.68	-0.98	25	0.334	-0.19
JGU	-0.26	0.53	-0.67	0.55	4.33	10	0.001**	1.20
WWU	-0.53	0.67	-1.11	0.56	4.82	25	5.90e-05***	0.92
NDSU	-0.33	0.64	-0.64	0.61	2.53	21	0.020	0.52
OSU	-0.64	0.46	-1.02	0.77	2.75	19	0.013*	0.59
UI	-0.35	0.73	-0.52	0.72	1.52	23	0.141	0.30
TUOS	-0.30	0.68	-0.58	0.74	2.12	14	0.052	0.52
UZH-MPR	-0.30	0.56	-0.85	0.70	6.85	29	1.58e-07***	1.22
UZH-NCN	-0.61	0.68	-1.14	0.85	4.54	19	2.24e-04***	0.97

12 Note. M = Mean. SD = Standard deviation. Df = Degrees of freedom.

13 \* $p < 0.02$ . \*\* $p < 0.01$ . \*\*\* $p < 0.001$

14

15 To assess the overall effect across labs, we performed a random-effects meta-analysis on the  
16 site-wise effect sizes (Figure 5). The meta-analysis revealed a robust set size effect (average  
17  $g_z = 0.63$ , 98%-CI = [0.25, 1.02],  $t = 4.63$ ,  $p = 0.001$ ), replicating the original effects. Thus,  
18 despite minor variability across individual labs, the meta-analytic evidence indicates a reliable  
19 increase in CDA negativity for set size 4 compared with set size 2.



1  
2 Figure 5. Direct Pipeline. Forest plot of the meta-analysis for set size 2 vs. 4 across laboratories. For  
3 each laboratory, the plot shows the effect size estimate Hedges'  $g$  with its standard error and  
4 corresponding 98% confidence interval. Squares represent lab-specific effect size estimates, with  
5 square size proportional to the inverse-variance weight, indicating the relative contribution of each  
6 laboratory to the pooled estimate (larger squares reflect greater weight due to higher precision).  
7 Horizontal lines denote 98% confidence intervals. The random-effects meta-analytic estimate (Hartung-  
8 Knapp adjustment) is shown as a diamond with its 98% confidence interval, together with the 98%  
9 prediction interval. Between-laboratory heterogeneity is quantified using  $I^2$ , with the associated  $p$ -value  
10 reported.

### 11 Hypothesis #1.2. CDA amplitude increase from set size 2 to 6

12 We next compared CDA amplitudes between set sizes 2 and 6 (Table 11). Post-hoc t-tests  
13 revealed statistically significant differences between set sizes 2 and 6 only in 3 out of 10 labs  
14 (WWU, and both UZH labs).  
15

16  
17 Table 11. Direct Pipeline. Post-hoc t-tests for the set size 2 vs. 6 comparison across laboratories.

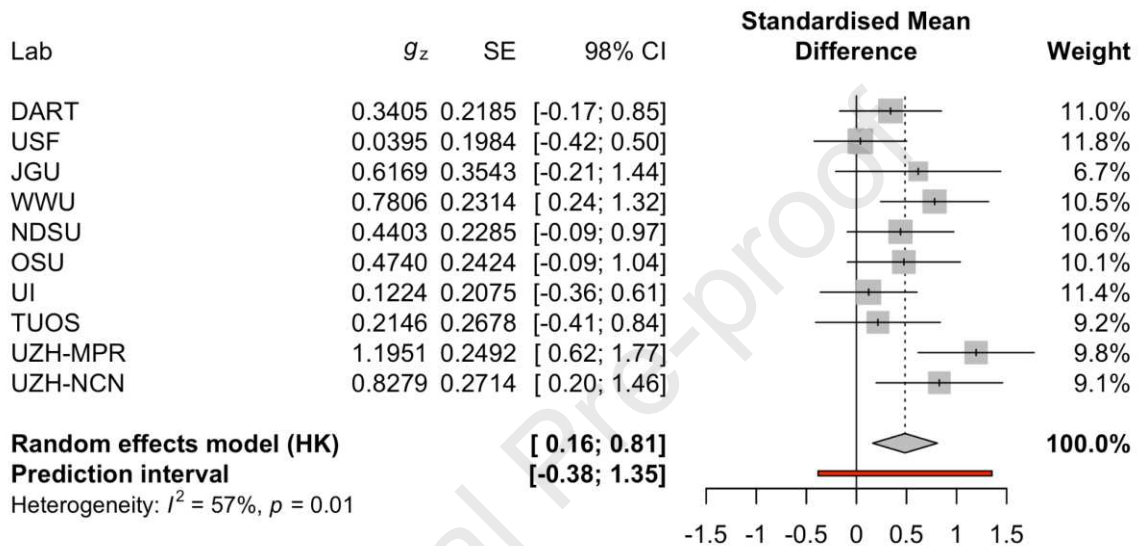
Lab	$M_{S_{z2}}$	$SD_{S_{z2}}$	$M_{S_{z6}}$	$SD_{S_{z6}}$	t-value	df	p-value	Hedges' $g_z$
DART	-0.61	0.62	-0.75	0.75	1.69	22	0.105	0.34
USF	-0.31	0.75	-0.35	0.71	0.21	25	0.837	0.04
JGU	-0.26	0.53	-0.53	0.72	2.22	10	0.051	0.62
WWU	-0.53	0.67	-1.05	0.69	4.10	25	3.78e-04***	0.78
NDSU	-0.33	0.64	-0.58	0.56	2.14	21	0.044	0.44
OSU	-0.64	0.46	-0.97	0.79	2.21	19	0.040	0.47
UI	-0.35	0.73	-0.44	0.68	0.62	23	0.541	0.12
TUOS	-0.30	0.68	-0.43	0.81	0.88	14	0.394	0.21
UZH-MPR	-0.30	0.56	-0.82	0.63	6.72	29	2.25e-07***	1.20
UZH-NCN	-0.61	0.68	-0.99	0.82	3.86	19	0.001**	0.83

18 Note. M = Mean. SD = Standard deviation. Df = Degrees of freedom.

1 \* $p < 0.02$ . \*\* $p < 0.01$ . \*\*\* $p < 0.001$

2

3 To quantify the overall effect across laboratories, we again performed a random-effects meta-  
 4 analysis on the site-level effect sizes (Figure 6). This analysis yielded a significant overall  
 5 difference between set sizes 2 and 6 (average  $g_z = 0.53$ , 98%-CI = [0.16, 0.81],  $t = 4.24$ ,  $p =$   
 6 0.002), replicating the original effect. These results indicate a reliable increase in CDA  
 7 negativity for set size 6 compared with set size 2 across labs.



8

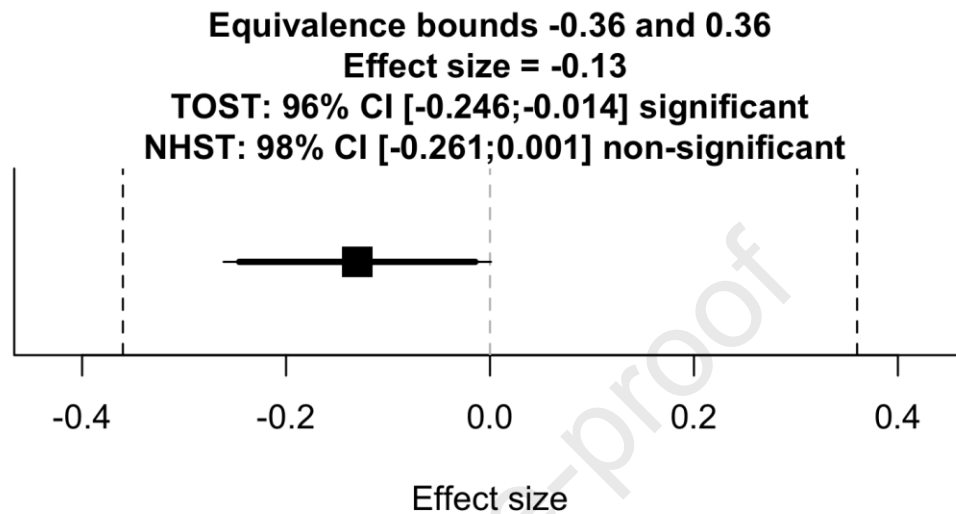
9 Figure 6. Direct Pipeline. Forest plot of the meta-analysis for set size 2 vs. 6 across laboratories. For  
 10 each laboratory, the plot shows the effect size estimate Hedges'  $g$  with its standard error and  
 11 corresponding 98% confidence interval. Squares represent lab-specific effect size estimates, with  
 12 square size proportional to the inverse-variance weight, indicating the relative contribution of each  
 13 laboratory to the pooled estimate (larger squares reflect greater weight due to higher precision).  
 14 Horizontal lines denote 98% confidence intervals. The random-effects meta-analytic estimate (Hartung-  
 15 Knapp adjustment) is shown as a diamond with its 98% confidence interval, together with the 98%  
 16 prediction interval. Between-laboratory heterogeneity is quantified using  $I^2$ , with the associated  $p$ -value  
 17 reported.

18

### 19 Hypothesis #1.3. Equivalence between CDA amplitude in set sizes 4 and 6

20 To assess whether CDA amplitudes for set sizes 4 and 6 were statistically equivalent, we  
 21 conducted an equivalence test on the random-effects meta-analytic effect size, using  
 22 equivalence bounds defined by the preregistered SESOI (Cohen's  $d = \pm 0.36$ ; Figure 7). The  
 23 TOST was performed on the pooled Hedges'  $g$  estimate across all labs and its standard error  
 24 (average  $g_z = -0.13$ , 98%-CI = [-0.29, 0.03],  $t = -2.31$ ,  $p = 0.046$ ), with significance evaluated  
 25 using 96% confidence intervals ( $\alpha = 0.02$ ), in line with the statistical framework adopted across  
 26 analyses. The TOST indicated statistical equivalence between set sizes 4 and 6 ( $Z = 4.06$ ,

1 96%-CI = [-0.25, -0.01],  $p = 2.50e-5$ ). The traditional null-hypothesis significance test did not  
 2 reach significance ( $Z = -2.31$ , 98%-CI = [-0.26, 0.003],  $p = 0.021$ ), indicating no detectable  
 3 difference between set sizes 4 and 6. Taken together, our analyses indicate that CDA  
 4 amplitudes for set sizes 4 and 6 do not differ significantly and fall within the preregistered  
 5 equivalence bounds, replicating the pattern reported in the original study.



6  
 7 Figure 7. Direct Pipeline. Equivalence test results for the meta-analytic comparison of set sizes 4 and  
 8 6 (H1.3). Mean difference (black square) is shown together with 96% TOST confidence intervals (thick  
 9 horizontal lines) and 98% NHST confidence intervals (thin horizontal lines). Equivalence bounds were  
 10 set at  $d = -0.36$  and  $d = 0.36$  based on the small-telescope approach (dashed vertical lines).

### 11 12 3. 4. 2. 2. Bayesian Approach

13  
 14 To complement the frequentist analyses, we estimated a Bayesian generalized linear mixed  
 15 model predicting CDA amplitude from set size (factor with 3 levels: set sizes 2, 4, and 6;  
 16 reference level: set size 2) while including subject and laboratory as random intercepts. Again,  
 17 the Bayesian linear mixed model additionally controlled for gender and handedness, in  
 18 accordance with Stage 1 Registered Report.

$$19 \quad \text{CDA} \sim \text{SetSize} * \text{Gender} * \text{Handedness} + (1|\text{Subject}) + (1|\text{Lab})$$

20  
 21  
 22 The posterior distributions provided clear evidence for set-size-dependent changes in CDA  
 23 amplitude. Relative to set size 2 (reference level), CDA amplitudes were more negative for  
 24 both set size 4 (Estimate = -0.22, 98%-CI = [-0.31, -0.13]) and set size 6 (Estimate = -0.19,  
 25 98%-CI = [-0.28, -0.10]). When re-leveling the model to use set size 4 as the reference, the  
 26 posterior estimate for the contrast between set sizes 6 and 4 was small (Estimate = 0.03, 98%-

1 CI = [-0.04, 0.11]), with the credible interval including zero, indicating no evidence for a  
2 difference in CDA amplitude between the two larger set sizes. None of the covariates showed  
3 credible associations with CDA amplitude: the 98%-CIs for gender and handedness predictors  
4 all included zero, indicating no evidence that these factors meaningfully contributed to  
5 variability in CDA amplitude. Taken together, the Bayesian results closely mirror the  
6 frequentist findings, providing convergent evidence that CDA amplitudes reliably increase  
7 from set size 2 to 4 and from set size 2 to 6 across laboratories.

8  
9 Finally, we implemented the Bayesian sequential updating procedure to evaluate the stability  
10 of the set-size effects as evidence accumulated across laboratories. For the first site, we fit a  
11 Bayesian regression model predicting CDA amplitude from set size using weakly informative  
12 (uniform) priors. The posterior from this model was then used as the prior for the next lab, and  
13 this process was repeated iteratively across all ten sites, allowing the posterior distribution to  
14 be continuously updated as new data were incorporated. The final posterior distribution after  
15 integrating evidence across all labs provided clear support for the predicted set-size effects.  
16 Relative to set size 2 (reference level), CDA amplitudes were credibly more negative for both  
17 set size 4 (Estimate = -0.24, 98%-CI = [-0.30, -0.18]) and set size 6 (Estimate = -0.18, 98%-  
18 CI = [-0.24, -0.13]). In both cases, the 98% credible intervals excluded zero, indicating strong  
19 cumulative evidence for load-dependent increases in CDA negativity as memory load  
20 increases. When re-leveling the model to use set size 4 as the reference, the posterior  
21 estimate for the contrast between set sizes 6 and 4 was small (Estimate = 0.06, 98%-CI = [-  
22 0.001, 0.12]), with the credible interval including zero, indicating no evidence for a difference  
23 in CDA amplitude between the set sizes 4 and 6.

### 24 25 3. 4. 3. Hypothesis #2: Correlation of CDA Increase with VWM Capacity

26  
27 In the following analyses, we investigated whether individual differences in behavioral  
28 performance in the change detection task are related to CDA amplitude increases from 2 to 4  
29 items and from 4 to 6 items, using both frequentist and Bayesian statistical models to test our  
30 preregistered hypotheses.

#### 31 32 3. 4. 3. 1. Frequentist Approach

##### 33 34 Hypothesis #2.1. Correlation between CDA increase from set size 2 to 4 and VWM capacity

35 First, we examined whether the increase in CDA amplitude from set size 2 to 4 correlates with  
36 VWM capacity, computed as average K-score over set sizes 2, 4, and 6. The correlation

1 coefficients were computed separately for each lab and then synthesized using a random-  
 2 effects meta-analysis. The correlations did not reach statistical significance in a single lab (all  
 3  $p > 0.02$ ). As shown in Table 12 and Figure 8, individual lab correlations varied in magnitude  
 4 and direction, with confidence intervals often spanning zero.

5

6 Table 12. Direct Pipeline. Correlations between CDA amplitude increase from set sizes 2 to 4 and VWM  
 7 capacity across laboratories

Lab	$M_{S22}$	$SD_{S22}$	$M_{S24}$	$SD_{S24}$	Pearson's r	p-value
DART	-0.25	0.36	2.22	0.52	0.38	0.074
USF	0.17	0.87	1.75	0.60	-0.33	0.094
JGU	-0.41	0.31	1.85	0.66	0.42	0.200
WWU	-0.59	0.62	2.69	0.45	0.14	0.485
NDSU	-0.31	0.58	1.67	0.61	0.16	0.488
OSU	-0.38	0.61	2.29	0.46	0.44	0.050
UI	-0.18	0.57	1.76	0.48	0.15	0.488
TUOS	-0.29	0.52	1.85	0.52	0.26	0.342
UZH-MPR	-0.55	0.44	2.26	0.62	-0.05	0.792
UZH-NCN	-0.54	0.53	2.36	0.61	0.04	0.862

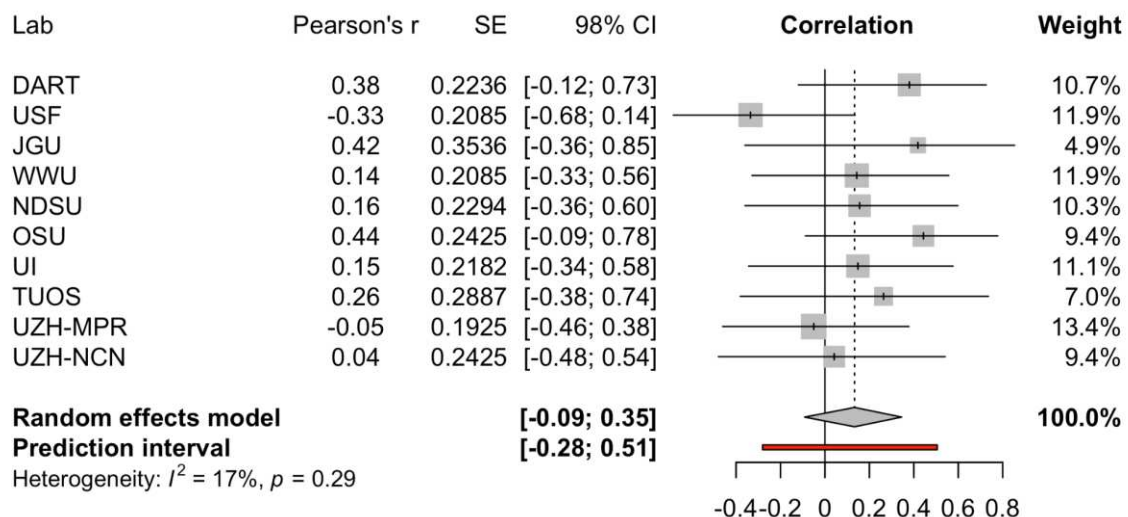
8 Note. M = Mean. SD = Standard deviation.

9 \* $p < 0.02$ . \*\* $p < 0.01$ . \*\*\* $p < 0.001$

10

11 Across the ten labs, the random-effects meta-analysis yielded a small positive correlation that  
 12 did not reach statistical significance ( $r = 0.13$ , 98%-CI = [-0.09, 0.36],  $t = 1.68$ ,  $p = 0.128$ ),  
 13 indicating no reliable evidence for a positive association between VWM capacity defined as a  
 14 K-score and the CDA increase from set size 2 to 4. Taken together, these results do not  
 15 provide evidence for the positive relationship predicted by Hypothesis 2.1, nor do they  
 16 replicate the strong association reported in the original study.

17



18

1 Figure 8. Direct Pipeline. Forest plot of the meta-analysis for the correlation between amplitude increase  
2 from set size 2 to 4 and VWM capacity defined as K-score. For each laboratory, the plot shows  
3 Pearson's  $r$  correlation coefficient with its standard error and corresponding 98% confidence interval.  
4 Squares represent lab-specific effect size estimates, with square size proportional to the inverse-  
5 variance weight, indicating the relative contribution of each laboratory to the pooled estimate (larger  
6 squares reflect greater weight due to higher precision). Horizontal lines denote 98% confidence  
7 intervals. The random-effects meta-analytic estimate (Hartung-Knapp adjustment) is shown as a  
8 diamond with its 98% confidence interval, together with the 98% prediction interval. Between-laboratory  
9 heterogeneity is quantified using  $I^2$ , with the associated p-value reported.

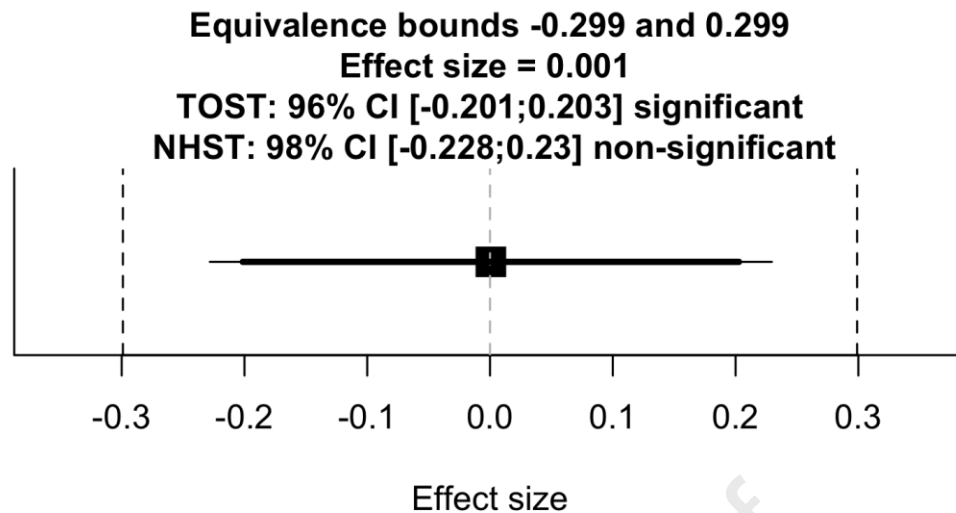
10  
11 As a sensitivity analysis, we conducted an analogous meta-analysis using  $d'$  instead of K-  
12 score to quantify performance. Across the ten laboratories, the random-effects meta-analysis  
13 revealed a small positive but non-significant correlation ( $r = 0.09$ , 98% CI = [-0.14, 0.32],  $t =$   
14  $1.10$ ,  $p = 0.301$ ), supporting the results of the analyses using K-score.

15  
16 Given the comparatively poor data quality observed for the USF laboratory (Table 7), we  
17 conducted an exploratory sensitivity analysis to assess whether this site disproportionately  
18 influenced the correlation between the CDA increase from set size 2 to 4 and individual VWM  
19 capacity. Specifically, we repeated the full set of analyses for H2.1 after excluding the USF  
20 dataset. Across the 9 remaining laboratories, the random-effects meta-analysis revealed a  
21 slightly stronger positive correlation between CDA increase and K score that reached  
22 statistical significance ( $r = 0.19$ , 98%-CI [0.02, 0.35],  $t = 3.16$ ,  $p = 0.013$ ). In contrast, the  
23 corresponding analysis using  $d'$  as the measure of VWM capacity showed a smaller, non-  
24 significant association ( $r = 0.13$ , 98%-CI [-0.10, 0.35],  $t = 1.63$ ,  $p = 0.142$ ).

#### 26 Hypothesis #2.2. Correlation between CDA increase from set size 4 to 6 and VWM capacity

27 Subsequently, to evaluate whether the association between VWM capacity and the CDA  
28 amplitude increase from set size 4 to 6 was small enough to be considered statistically  
29 equivalent, we conducted an equivalence test on the random-effects meta-analytic correlation,  
30 using the preregistered SESOI ( $r = \pm 0.29$ ) derived from the small telescopes approach (Figure  
31 9). The TOST indicated statistical equivalence ( $Z = -3.03$ , 96%-CI = [-0.20, 0.20],  $p = 0.002$ ).  
32 The traditional null-hypothesis significance test did not reach significance ( $Z = 0.01$ , 98%-CI =  
33 [-0.23, 0.23],  $p = 0.994$ ), indicating no detectable correlation between VWM capacity and the  
34 CDA increase from set sizes 4 to 6. Taken together, these results show that the CDA increase  
35 from 4 to 6 items is both statistically indistinguishable from zero and statistically equivalent  
36 within the preregistered bounds, thereby replicating the pattern reported in the original study.

37



1

2 Figure 9. Direct Pipeline. Equivalence test for the correlation between VWM capacity and CDA  
 3 increase from set size 4 to 6 (H2.2). Mean difference (black square) is shown together with 96%  
 4 TOST confidence intervals (thick horizontal lines) and 98% NHST confidence intervals (thin horizontal  
 5 lines). Equivalence bounds were set at  $d = -0.299$  and  $d = 0.299$  based on the small-telescope  
 6 approach (dashed vertical lines).

7

### 8 3. 4. 3. 2. Bayesian Approach

9

#### 10 Hypothesis #2.1. Correlation between CDA increase from set size 2 to 4 and VWM capacity

11 To complement the frequentist analyses, we examined the relationship between VWM  
 12 capacity and the CDA increase from set size 2 to 4 within a Bayesian framework. We fitted a  
 13 Bayesian linear mixed model with VWM capacity as the predictor and the CDA difference (set  
 14 size 4 minus set size 2) as the outcome, including laboratory as a random intercept, and  
 15 gender and handedness as covariates. Consistent with our preregistration, the model used a  
 16 prior centered on the correlation reported in the original study ( $r = 0.78$ ).

17

$$18 \quad \text{CDA}_{4-2} \sim K * \text{Gender} * \text{Handedness} + (1|\text{Lab})$$

19

20 The posterior distribution provided no evidence for a relationship between the CDA increase  
 21 and VWM capacity (Estimate = 0.16, 98%-CI = [-0.06, 0.38]). The posterior mean was close  
 22 to zero, and the credible interval spanned both positive and negative values, indicating a high  
 23 probability that any true association is negligible. The posterior distributions for all covariates  
 24 (gender and handedness) likewise overlapped zero, suggesting that none of the covariates  
 25 contributed meaningfully to the model.

26

1 Finally, we applied the Bayesian sequential updating procedure to evaluate whether evidence  
 2 for a relationship between VWM capacity and the CDA increase from set size 2 to 4 would  
 3 accumulate across laboratories. Starting with weakly informative priors for the first dataset,  
 4 the posterior from each lab was used as the prior for the subsequent lab, allowing the evidence  
 5 to be integrated iteratively across all ten sites. The final posterior distribution converged tightly  
 6 around zero, with the posterior mean close to zero and the 98% credible interval including  
 7 both positive and negative values (Estimate = 0.14, 98%-CI [-0.03, 0.31]). Taken together, the  
 8 sequential updating analysis confirms the results from the frequentist, and standard Bayesian  
 9 models: across accumulating evidence from all laboratories, there is no indication that CDA  
 10 increase from set size 2 to 4 is related to VWM capacity.

11

### 12 Hypothesis #2.2. Correlation between CDA increase from set size 4 to 6 and VWM capacity

13 Finally, we examined the relationship between VWM capacity and the CDA increase from set  
 14 size 4 to 6 within a Bayesian framework. We fitted a Bayesian linear mixed model with VWM  
 15 capacity as the predictor and the CDA difference (set size 6 minus set size 4) as the outcome,  
 16 including laboratory as a random intercept. The original study did not report the correlation  
 17 between amplitude increase from set size 4 to 6, therefore, we used weakly informative  
 18 (uniform) priors.

19

$$20 \quad \text{CDA}_{6-4} \sim K * \text{Gender} * \text{Handedness} + (1|\text{Lab})$$

21

22 The posterior distribution provided no evidence for a relationship between the CDA increase  
 23 from set size 4 to 6 and VWM capacity (Estimate = 0.04, 98%-CI = [-0.17, 0.26]). The posterior  
 24 mean was close to zero, and the credible interval spanned both positive and negative values,  
 25 indicating a high probability that any true association is negligible. The posterior distributions  
 26 for all covariates (gender and handedness) likewise overlapped zero, suggesting that none of  
 27 the covariates contributed meaningfully to the model.

28

29 Consistent with this result, the Bayesian sequential updating analysis yielded an estimate  
 30 close to zero (Estimate = 0.002, 98%-CI = [-0.17, 0.18]), with the credible interval again  
 31 spanning both positive and negative values, providing no evidence for a reliable association  
 32 between the CDA increase from set size 4 to 6 and VWM capacity.

## 1 Discussion

2 In this multi-site replication study, we aimed to evaluate two central claims from Vogel and  
3 Machizawa (2004): that the CDA reliably tracks increases in VWM load, and that the  
4 magnitude of the CDA increase from set size 2 to 4 correlates with individual VWM capacity.  
5 Consistent with the original findings, we reproduced a clear CDA lateralization effect and  
6 robust set size increases in CDA amplitude from set size 2 to 4 and from set size 2 to 6,  
7 demonstrating that the component is a robust neural marker of VWM maintenance across  
8 diverse samples, EEG systems, and testing environments. In contrast, we found no compelling  
9 evidence for the predicted positive correlation between CDA amplitude increase from set size  
10 2 to 4 and VWM capacity. With the exception of a small effect that reached statistical credibility  
11 only in the Bayesian sequential-updating analysis of the advanced pipeline, all other estimates  
12 were smaller than those reported in the original study and previous meta-analyses. Taken  
13 together, these findings indicate that while the CDA is a robust indicator of memory load, its  
14 value as a reliable marker of individual differences in VWM capacity is likely substantially more  
15 limited than previously assumed.

16

### 17 CDA as a Marker of Memory Load (Hypothesis #1)

18

19 Across laboratories, our results provide strong support for the robustness of the CDA as a  
20 neural marker of visual working memory load, in line with a large body of prior work (Adam et  
21 al., 2018; Diamantopoulou et al., 2011; Drew & Vogel, 2008; Feldmann-Wüstefeld, 2021;  
22 Feldmann-Wüstefeld et al., 2018; Hakim et al., 2019, 2020; Kang & Woodman, 2014; Kundu  
23 et al., 2013; Kuo et al., 2012; Lefebvre et al., 2013; Leonard et al., 2013; Luria et al., 2016;  
24 Ngiam et al., 2021; Roy & Faubert, 2023; Störmer et al., 2013; Tröndle & Langer, 2024;  
25 Tsubomi et al., 2013; Unsworth et al., 2015; Villena-González et al., 2020; Vogel &  
26 Machizawa, 2004). In both the direct and advanced pipelines, all (or all but one) laboratories  
27 showed a robust contralateral-ipsilateral asymmetry, and the meta-analytic effect sizes were  
28 large ( $g_z \approx -0.95$ ), confirming that a sustained contralateral negativity emerges consistently  
29 during the retention interval. Importantly, CDA amplitudes increased systematically from set  
30 size 2 to 4 and from 2 to 6, with random-effects meta-analyses yielding medium-to-large  
31 effects ( $g_z \approx 0.45-0.65$ ) despite some variability in statistical significance at the level of  
32 individual sites, which is expected when the same underlying effect is tested repeatedly across  
33 multiple independent samples. Equivalence tests further indicated that CDA amplitudes for set  
34 sizes 4 and 6 were statistically equivalent and fell within the preregistered equivalence bounds,  
35 replicating the plateau pattern reported by Vogel and Machizawa (2004), which parallels the  
36 average human working memory capacity limit of around four items. The Bayesian mixed

1 models and sequential updating analyses converged with the frequentist results, providing  
2 clear evidence for load-dependent increases from set size 2 to 4 and 2 to 6, but no credible  
3 difference between 4 and 6 items. Together with the high dependability coefficients and low  
4 SME values observed at most sites, these findings demonstrate that the CDA can be  
5 measured with good precision and internal consistency across different samples, EEG  
6 systems, laboratories, and recordings environments. Thus, at the process level, our multi-site  
7 study confirms that the CDA is a robust and psychometrically reliable electrophysiological  
8 index of the number of items held in visual working memory.

## 9 10 Correlation of CDA Increase with VWM Capacity (Hypothesis #2)

11  
12 In contrast to the robust set-size effects, our multi-lab replication provided only limited  
13 evidence for the predicted positive association between VWM capacity and the CDA increase  
14 from set size 2 to 4. In both, direct and advanced pipelines, correlations were non-significant  
15 in every laboratory, and the random-effects meta-analysis yielded a small, non-significant  
16 estimate. The Bayesian mixed model likewise produced a posterior centered near zero,  
17 indicating that the current data provided little evidence for a positive relationship between CDA  
18 amplitude increase and VWM capacity. In the ICA pipeline, frequentist correlations were again  
19 non-significant across all sites. However, the Bayesian mixed model and the sequential  
20 Bayesian updating both indicated a small positive association, with an estimated effect of  
21 approximately  $r \approx 0.20$ . Notably, an effect of  $r \approx 0.20$  is small, accounting for roughly 4% of the  
22 variance, and thus represents only a weak association between neural and behavioral  
23 measures. Taken together, results from all pipelines converge on the conclusion that the CDA  
24 increase does not reliably track individual differences in VWM capacity and that any true  
25 relationship, if present, is substantially smaller and less robust than previously suggested.

26  
27 The lack of a strong correlation between CDA and behavior deviates from a major finding of  
28 the original (2004) study and from later meta-analyses (Luria et al., 2016; Roy & Faubert,  
29 2023). For instance, the meta-analysis of 11 studies by Luria et al. (2016) yielded a combined  
30 correlation of  $r = 0.596$  (98%-CI = [0.51, 0.67]), suggesting a robust relationship between CDA  
31 amplitude increase and VWM capacity in the existing literature. Similarly, Roy and Faubert  
32 (2023) reported three medium-to-strong correlations ranging from  $r = 0.26$  to 0.45.

33  
34 One likely source of discrepancy between our findings and those in the existing literature is  
35 the limited sample size of most published studies examining the correlations between CDA  
36 and VWM capacity, together with substantial between-study variability in how the relationship

1 between CDA and VWM capacity is operationalized. Previous studies have correlated VWM  
2 capacity with differences between two different set sizes (e.g., 4-2, 3-1) or with CDA amplitude  
3 at a single set size (e.g., set size 3 or 6). Psychometrically, the reliability of a difference score  
4 is inherently limited when the two contributing variables are highly correlated, as is likely for  
5 capacity estimates at set sizes 4 and 6 provided that variability exists at both levels.  
6 Furthermore, difference scores may capture a distinct, and often smaller, portion of variance  
7 in individual differences in VWM than performance assessed at a single set size or across  
8 conditions, because shared variance between conditions is effectively removed. Given that  
9 mean capacity typically asymptotes around set size 4, the difference between set sizes 4 and  
10 6 may, for many individuals, reflect compensatory strategies or the ability to maintain  
11 information beyond nominal capacity limits rather than core storage capacity per se. These  
12 analytic choices complicate direct comparisons across studies and may contribute to variability  
13 in reported effect sizes.

14  
15 Additionally, the majority of studies included in Luria et al.'s (2016) meta-analysis were small  
16 (e.g., Diamantopoulou et al., N = 14; Drew et al., N = 33/18; Jost et al., N = 25; Kang &  
17 Woodman, N = 24; Kundu et al., N = 30; Kuo et al., N = 18; Lefebvre et al., N = 39; Leonard  
18 et al., N = 23; Störmer et al., N = 35; Tsubomi et al., N = 25), with sample sizes between 14  
19 and 39 participants (Table 13). Such sample sizes fall far below the 250 participants needed  
20 for correlation estimates to stabilize (Schönbrodt & Perugini, 2013), making inflated and  
21 unstable correlations statistically likely. Consistent with this interpretation, the few larger  
22 studies in the literature report substantially smaller effects. For instance, Unsworth et al. (2015;  
23 N = 170) observed  $r = 0.33$  between the CDA amplitude at set size 6 (i.e., instead of a  
24 difference between 2 set sizes) and VWM capacity, Adam et al. (2018; N = 72) reported  $r =$   
25  $0.27$  (again, CDA amplitude at set size 6) and Tröndle and Langer (2024; N = 55) found only  
26 a modest association ( $r = 0.22$ ; CDA amplitude at set size 3). These correlations are closer  
27 to our own meta-analytic estimate and suggest that the true relationship, if present, is likely  
28 smaller, around  $r = 0.20-0.30$ . For context, when directly using CDA amplitude at individual  
29 set sizes instead of the CDA increase from set size 2 to 4, the meta-analytic correlations in  
30 our study were small and showed an inconsistent pattern across preprocessing pipelines. At  
31 set size 2, the meta-analytic correlations were near zero and non-significant (direct:  $r = 0.09$ ,  
32  $p = 0.234$ ; advanced:  $r = 0.06$ ,  $p = 0.382$ ; ICA:  $r = 0.02$ ,  $p = 0.786$ ). At set size 4, correlations  
33 were slightly larger but remained non-significant (direct:  $r = 0.17$ ,  $p = 0.079$ ; advanced:  $r =$   
34  $0.16$ ,  $p = 0.122$ ; ICA:  $r = 0.17$ ,  $p = 0.055$ ). The strongest associations were observed at set  
35 size 6, where correlations were small but statistically significant across pipelines (direct:  $r =$   
36  $0.17$ ,  $p = 0.019$ ; advanced:  $r = 0.15$ ,  $p = 0.034$ ; ICA:  $r = 0.15$ ,  $p = 0.022$ ). Overall, this pattern

1 further supports the conclusion that the association between CDA amplitude and VWM  
2 capacity is modest at best.

3

4 More broadly, the multitude of analytical choices researchers face when quantifying the  
5 relationships between CDA and behavior is referred to as researcher degrees of freedom  
6 (Gelman & Loken, 2023; Simmons et al., 2011; Trübtschek et al., 2023). When researchers'  
7 decisions are made post-hoc or without explicit justification, they can inflate variability, obscure  
8 true effect magnitudes, and contribute to inconsistent findings across studies (Götz et al.,  
9 2024; Simmons et al., 2011). Standardizing CDA operationalizations, or systematically  
10 evaluating alternative definitions within, for example, a multiverse framework, would therefore  
11 be an important step toward determining whether, and under which conditions, the CDA  
12 indexes individual differences in VWM capacity (Götz et al., 2024; Sarma et al., 2024).

13

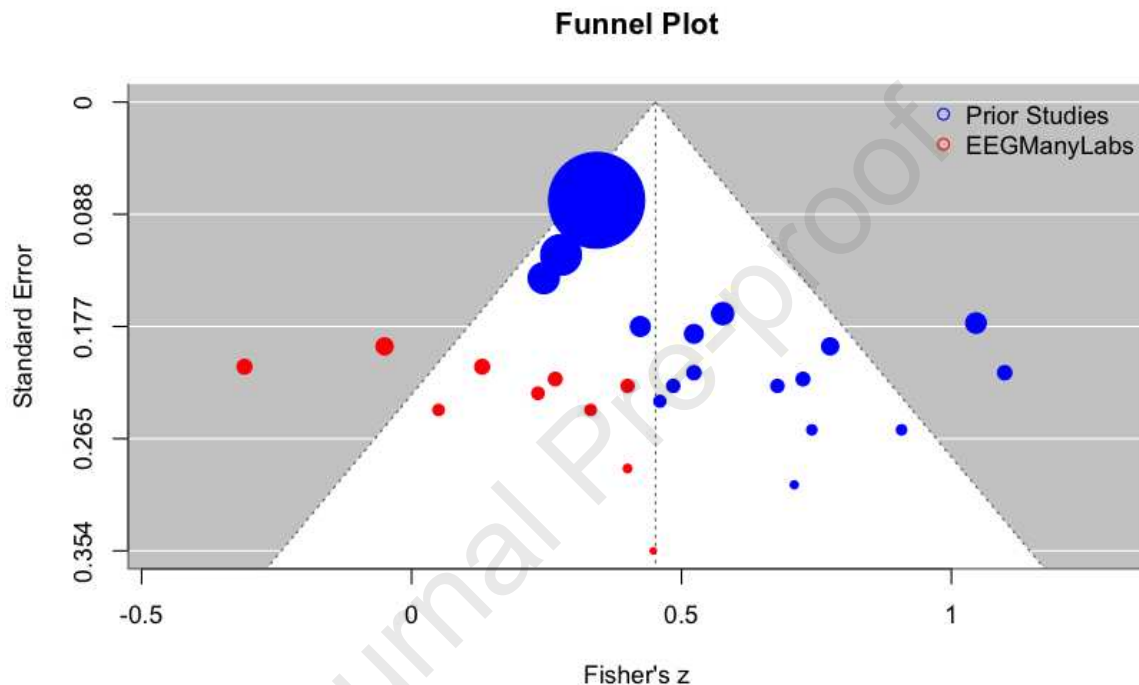
14 Another potential source of the discrepancy between our findings and previously reported  
15 medium-to-large effects is publication bias. We therefore conducted funnel-plot diagnostics  
16 (Sterne & Egger, 2001; Viechtbauer, 2010) on the complete set of 27 studies known to us that  
17 investigated the association between CDA and VWM capacity (Adam et al., 2018;  
18 Diamantopoulou et al., 2011; Drew & Vogel, 2008; Feldmann-Wüstefeld, 2021; Hakim et al.,  
19 2019, 2020; Kang & Woodman, 2014; Kundu et al., 2013; Kuo et al., 2012; Lefebvre et al.,  
20 2013; Leonard et al., 2013; Luria et al., 2016; Ngiam et al., 2021; Roy & Faubert, 2023;  
21 Störmer et al., 2013; Tröndle & Langer, 2024; Tsubomi et al., 2013; Unsworth et al., 2015;  
22 Villena-González et al., 2020; Vogel & Machizawa, 2004). For an overview of the studies  
23 published prior to EEGManyLabs, please refer to Table 13.

24

25 First, when restricting the meta-analysis to all prior #EEGManyLabs studies ( $k = 17$ ; blue  
26 studies on Figure 10), the random-effects model yielded a large, pooled effect (Fisher's  $z =$   
27  $0.58$ , 95%-CI =  $[0.45, 0.72]$ ,  $p = 2.22e-18$ ), accompanied by substantial heterogeneity ( $I^2 =$   
28  $56.73\%$ ). Funnel-plot asymmetry was pronounced in this restricted dataset with the Egger  
29 regression test statistically significant ( $z = 2.83$ ,  $p = 0.005$ ), providing strong evidence for  
30 small-study effects consistent with selective reporting. In contrast, when the full dataset  
31 including the #EEGManyLabs replication was analyzed ( $k = 27$ ), evidence for publication bias  
32 was strongly attenuated. The Egger regression test was non-significant ( $z = 0.72$ ,  $p = 0.474$ ),  
33 and the limit estimate as the standard error approached zero was  $b = 0.28$  (98%-CI =  $[-0.19,$   
34  $0.76]$ ), indicating no reliable inflation of effects in smaller samples. The random-effects model  
35 yielded a pooled Fisher's  $z$  of  $0.45$  (95%-CI =  $[0.33, 0.57]$ ,  $p = 4.64e-13$ ), corresponding to a  
36 significant medium effect. However, heterogeneity was substantial ( $I^2 = 63.43\%$ ), indicating  
37 considerable variability in effect sizes across studies.

1  
2  
3  
4  
5  
6  
7

Taken together, these results indicate that the large correlations between VWM capacity and CDA amplitude reported in the early literature are substantially inflated by small-study effects, whereas the inclusion of the present large-scale multi-laboratory data strongly attenuates both effect size estimates and evidence for publication bias. These results highlight the need for large-scale, well-powered, and preregistered studies to approximate the true magnitude and robustness of a given effect.



8  
9  
10  
11  
12  
13  
14  
15

Figure 10. Funnel plot of the CDA-VWM capacity correlations across all 27 studies. Each point represents an individual study, plotted as Fisher's  $z$  against its standard error, and the size of each point is proportional to the study's sample size. The dashed vertical line indicates the pooled random-effects estimate (Fisher's  $z = 0.45$ ), and diagonal dashed lines denote the 98% pseudo-confidence limits around the summary effect. Correlations from the EEGManyLabs datasets were derived from the direct replication.

## 16 Theoretical Implications for Interpreting the CDA

17  
18  
19  
20  
21  
22  
23

The dissociation between robust process-level effects and weak individual-differences correlations has important implications for the interpretation of CDA as a neural marker of VWM capacity. First, if the true brain-behavior correlation is modest ( $r \approx 0.20-0.30$  at best), then CDA measurements may not reliably capture the same sources of individual variation that drive behavioral performance differences. Therefore, researchers should exercise caution when interpreting individual differences in CDA as reflecting differences in VWM capacity, and

1 that CDA-based inferences about capacity changes may need to be reconsidered or at  
2 minimum supported by parallel behavioral evidence.

3

4 Second, the present findings highlight a fundamental theoretical dissociation: although CDA  
5 is highly sensitive to within-individual variations in memory load, it shows little consistent  
6 association with between-person differences in capacity. Several non-mutually exclusive  
7 explanations warrant consideration. One possibility is that CDA primarily reflects the  
8 maintenance or active storage component of VWM, whereas behavioral capacity estimates  
9 such as Cowan's K conflate maintenance with a host of other processes including encoding  
10 efficiency, attentional filtering during stimulus presentation, resistance to interference, and  
11 decision or comparison processes during retrieval. If individuals vary substantially in these  
12 non-maintenance components, then CDA amplitude might correlate only weakly with overall  
13 task performance despite being a valid index of the number of items actively maintained.  
14 Another related possibility is that individuals with lower CDA amplitudes may engage  
15 compensatory strategies that preserve behavioral performance. For instance, some  
16 individuals might rely more heavily on verbal recoding, hierarchical chunking, or a quality-  
17 quantity tradeoff (maintaining fewer items at higher precision) to achieve similar behavioral  
18 capacity estimates despite differences in the neural signature of visual maintenance indexed  
19 by CDA. Such strategic flexibility could obscure the brain-behavior relationship at the  
20 individual-differences level while leaving within-person load effects intact.

21

22 Third, it is important to recognize that between-person variance in CDA amplitude likely  
23 reflects a mixture of cognitive and non-cognitive sources. Anatomical and biophysical factors  
24 such as skull thickness, cortical geometry, sulcal depth, and electrode-to-source distance can  
25 all influence EEG amplitude measures and vary considerably across individuals. Critically,  
26 these anatomical factors contribute to stable between-person differences in CDA amplitude  
27 but are entirely orthogonal to cognitive capacity. Within-person experimental manipulations  
28 (such as varying set size) effectively control for these individual anatomical differences, which  
29 may explain why load effects are robust and replicable while individual-differences correlations  
30 are weak. Future work employing source localization methods or individual structural MRI data  
31 might help to disentangle cognitive from anatomical sources of CDA variance.

32

### 33 Conclusion

34

35 In summary, our large-scale, preregistered multi-lab replication provides evidence that the  
36 CDA is a reliable neural marker of visual working memory load, consistently reproducing the

1 contralateral-ipsilateral asymmetry and the characteristic load-dependent increases observed  
 2 in the seminal work. At the same time, our results do not support the widely cited claim that  
 3 individual differences in CDA amplitude reflect individual differences in visual working memory  
 4 capacity. Across pipelines, statistical frameworks, and laboratories, the association between  
 5 CDA increase and behavioral capacity was small, inconsistent, and far weaker than previously  
 6 reported in the literature. Together with evidence of substantial heterogeneity among past  
 7 studies, the present findings suggest that earlier reports of medium-to-large correlations were  
 8 likely inflated by small sample sizes, analytic flexibility, and variation in measurement  
 9 reliability. These results highlight the critical role of large, collaborative, preregistered projects  
 10 for resolving long-standing debates and establishing the true size and robustness of  
 11 foundational cognitive neuroscience effects, as well as the necessity of adequately powered  
 12 studies to obtain stable and interpretable estimates of individual differences.

13

14 Table 13. Overview of studies published prior to EEGManyLabs examining contralateral delay activity  
 15 and working memory capacity.

Study	N	K-score (Mean, SD)						CDA definition	r
		SZ1	SZ2	SZ3	SZ4	SZ6	Average K		
Adam et al., 2018, Exp1	72	0.95 (0.04)	—	2.41 (0.33)	—	2.53 (0.53)	2.62 (1.00)	SZ6	0.26
Diamantopoulou et al., 2011	14	NR	NR	NR	NR	—	1.73	SZ3 - SZ2	0.61
Drew and Vogel, 2008, Exp 4	33	NR	—	NR	—	—	NR	SZ3 - SZ1	0.48
Drew and Vogel, 2008, Exp 5	18	NR	—	NR	—	—	NR	SZ3 - SZ1	0.72
Feldmann-Wüstefeld et al., 2021	21	—	NR	—	NR	NR	2.5	Mean of SZ2 and SZ4; SZ4 - SZ2	0.43; 0.39
Jost et al., 2011 (Young)	25	NR	—	2.14	—	—	NR	NR	0.48
Kang and Woodman, 2014	24	NR	NR	—	NR	NR	2.59	SZ4 - SZ1	0.62
Kundu et al., 2013	30	—	NR	—	NR	—	2.16	NR	0.62
Kuo et al., 2012, Exp1	18	—	1.53 (0.16)	—	2.66 (0.58)	—	NR	SZ4 - SZ2	0.63
Lefebvre et al., 2013	39	—	1.9	—	3.3	3.75	NR	SZ6 - SZ2	0.52
Leonard et al., 2013	23	NR	—	NR	—	—	NR	SZ3 - SZ1	0.59
Stormer et al., 2013 (Young)	35	NR	—	2.1	—	—	NR	SZ3 - SZ1	0.40
Tröndle & Langer, 2024 (Young)	55	0.96 (0.05)	—	2.26 (0.42)	—	—	—	SZ3	0.24
Tsubomi et al., 2013	25	—	NR	—	NR	NR	NR	SZ4 - SZ2	0.64

Unsworth et al., 2015	170	—	NR	—	NR	1.90 (0.77)	NR	SZ6	0.30
Villena-Gonzalez et al., 2020	23	NR	NR	—	NR	—	2.3 (0.8)	SZ4 - SZ2	0.45
Vogel and Machizawa, 2004	36	NR	NR	NR	NR	NR	2.8	SZ4 - SZ2	0.78

1 Note. SZ = set size. NR = not reported. R = Pearson's correlation coefficient. SD = standard deviation.  
 2 SD values are shown only when they were reported in the publication. Often, publications do not report  
 3 K scores separately for each investigated set size, which is indicated as NR in the table. An em dash  
 4 (—) indicates that a specific set size was not investigated in the given study.

5

6

## 7 Acknowledgements

8 #EEGManyLabs is supported by a DFG (PA 4005/1-1) grant to YGP and a UKRI BBSRC grant  
 9 (BB/X008428/1) awarded to FM. The Lead Author, NL, is funded by the Swiss National  
 10 Science Foundation (SNSF) Grant 100014\_175875. HMS, HD, and AL are supported by the  
 11 Icelandic Research Fund (228916-051) and the University of Iceland Research Fund. CCvB  
 12 and SJ are supported by the Economic and Social Research Council (UK; ES/V013610/1).  
 13 MGM is supported by the Moonshot R&D Goal 9 (JPMJMS2296) and JST COI (JPMJCE1311,  
 14 JPMJCA2208). PEC is funded by a grant from the National Institute of Mental Health  
 15 (MH128208). These funding sources were not involved in data analysis, interpretation, or the  
 16 preparation of this submission for publication.

17

## 18 Competing Interests

19 The authors disclose no conflicts of interest related to this manuscript.

20

## 21 Author Contributions

22 Conceptualization: D.S., F.M., Y.G.P., M.G.M., W.X.Q.N., E.K.V., and N.L.

23 Data curation: D.S. and N.L.

24 Formal analysis: D.S. and N.L.

25 Funding acquisition: P.E.C., H.M.S., F.M., Y.G.P., V.S.S., A.-L.S., J.S.J., J.D.G., C.C.v.B.,  
 26 N.A.B., and N.L.

27 Investigation: D.S., P.E.C., H.M.S., H.D., A.L., H.A.R., Y.H.C., K.M.O., V.S.S., J.C.G.A., C.L.,  
 28 A.-L.S., A.-L.B., S.A.B., E.M.J., J.S.J., Z.L., Y.M.C., E.L., J.D.G., S.J., M.J., E.M., C.C.v.B.,  
 29 N.A.B., C.P., L.B., Y.H., W.X.Q.N., and N.L.

30 Methodology: F.M., Y.G.P., K.M.O., M.G.M., W.X.Q.N., E.K.V., and N.L.

31 Project administration: F.M., Y.G.P., and N.L.

- 1 Resources: D.S., P.E.C., H.M.S., F.M., Y.G.P., V.S.S., A.-L.S., J.S.J., J.D.G., C.C.v.B.,
- 2 N.A.B., W.X.Q.N., and N.L.
- 3 Software: D.S., W.X.Q.N., and N.L.
- 4 Supervision: P.E.C., H.M.S., F.M., Y.G.P., V.S.S., A.-L.S., J.S.J., J.D.G., C.C.v.B., N.A.B.,
- 5 and N.L.
- 6 Validation: D.S., F.M., Y.G.P., and N.L.
- 7 Visualization: D.S., H.M.S., and N.L.
- 8 Writing - original draft: D.S. and N.L.
- 9 Writing - review & editing: D.S., P.E.C., H.M.S., F.M., Y.G.P., H.D., A.L., H.A.R., Y.H.C.,
- 10 K.M.O., V.S.S., J.C.G.A., C.L., A.-L.S., A.-L.B., S.A.B., E.M.J., J.S.J., Z.L., Y.M.C., E.L.,
- 11 J.D.G., S.J., M.J., E.M., C.C.v.B., N.A.B., C.P., L.B., Y.H., M.G.M., W.X.Q.N., E.K.V., and N.L.
- 12
- 13

## 1 References

- 2 Adam, K. C. S., Robison, M. K., & Vogel, E. K. (2018). Contralateral delay activity tracks  
3 fluctuations in working memory performance. *Journal of Cognitive Neuroscience*, *30*(9),  
4 1229–1240. [https://doi.org/10.1162/jocn\\_a\\_01233](https://doi.org/10.1162/jocn_a_01233)
- 5 Asp, I. E., Störmer, V. S., & Brady, T. F. (2021). Greater visual working memory capacity for  
6 visually matched stimuli when they are perceived as meaningful. *Journal of Cognitive  
7 Neuroscience*, 1–17. [https://doi.org/10.1162/jocn\\_a\\_01693](https://doi.org/10.1162/jocn_a_01693)
- 8 Atkinson, R. C., & Shiffrin, R. M. (1968). Human memory: A proposed system and its control  
9 processes. In *Psychology of Learning and Motivation* (pp. 89–195). Elsevier.
- 10 Baddeley, A. (2003). Working memory: looking back and looking forward. *Nature Reviews.  
11 Neuroscience*, *4*(10), 829–839.
- 12 Baddeley, A. D. (1986). *Working Memory*. OUP Australia and.
- 13 Baddeley, A. D., Gathercole, S. E., & Papagno, C. (2017). The phonological loop as a  
14 language learning device. In *Exploring Working Memory* (pp. 164–198). Routledge.
- 15 Baddeley, A. D., & Hitch, G. (1974). Working Memory. In *Psychology of Learning and  
16 Motivation* (pp. 47–89). Elsevier.
- 17 Baddeley, A., Papagno, C., & Vallar, G. (1988). When long-term learning depends on short-  
18 term storage. *Journal of Memory and Language*, *27*(5), 586–595.
- 19 Baldwin, S. A., Larson, M. J., & Clayson, P. E. (2015). The dependability of  
20 electrophysiological measurements of performance monitoring in a clinical sample: A  
21 generalizability and decision analysis of the ERN and Pe: EEG dependability.  
22 *Psychophysiology*, *52*(6), 790–800. <https://doi.org/10.1111/psyp.12401>
- 23 Brady, T. F., Robinson, M. M., Williams, J. R., & Wixted, J. T. (2023). Measuring memory is  
24 harder than you think: How to avoid problematic measurement practices in memory  
25 research. *Psychonomic Bulletin & Review*, *30*(2), 421–449.  
26 <https://doi.org/10.3758/s13423-022-02179-w>
- 27 Brady, T. F., Störmer, V. S., & Alvarez, G. A. (2016). Working memory is not fixed-capacity:

- 1 More active storage capacity for real-world objects than for simple stimuli. *Proceedings*  
2 *of the National Academy of Sciences of the United States of America*, 113(27), 7459–  
3 7464. <https://doi.org/10.1073/pnas.1520027113>
- 4 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, 10(4), 433–436.  
5 <https://doi.org/10.1163/156856897x00357>
- 6 Brennan, R. L. (1992). Generalizability theory. *Educational Measurement Issues and*  
7 *Practice*, 11(4), 27–34. <https://doi.org/10.1111/j.1745-3992.1992.tb00260.x>
- 8 Brisson, B., & Jolicoeur, P. (2007). A psychological refractory period in access to visual  
9 short-term memory and the deployment of visual-spatial attention: multitasking  
10 processing deficits revealed by event-related potentials. *Psychophysiology*, 44(2), 323–  
11 333. <https://doi.org/10.1111/j.1469-8986.2007.00503.x>
- 12 Bürkner, P.-C. (2017). Brms: An R package for Bayesian multilevel models using Stan.  
13 *Journal of Statistical Software*, 80(1). <https://doi.org/10.18637/jss.v080.i01>
- 14 Caldwell, A. R. (2022). Exploring equivalence testing with the updated TOSTER R package.  
15 In *PsyArXiv*. <https://doi.org/10.31234/osf.io/ty8de>
- 16 Caplan, D., & Waters, G. S. (1999). Verbal working memory and sentence comprehension.  
17 *The Behavioral and Brain Sciences*, 22(1), 77–94; discussion 95–126.
- 18 Champely, S. (2020). *Basic Functions for Power Analysis [R package pwr version 1.3-0]*.  
19 <https://CRAN.R-project.org/package=pwr>
- 20 Clayson, P. E., Baldwin, S. A., & Larson, M. J. (2021). Evaluating the internal consistency of  
21 subtraction-based and residualized difference scores: Considerations for psychometric  
22 reliability analyses of event-related potentials. *Psychophysiology*, 58(4), e13762.  
23 <https://doi.org/10.1111/psyp.13762>
- 24 Clayson, P. E., Brush, C. J., & Hajcak, G. (2021). Data quality and reliability metrics for  
25 event-related potentials (ERPs): The utility of subject-level reliability. *International*  
26 *Journal of Psychophysiology: Official Journal of the International Organization of*  
27 *Psychophysiology*, 165, 121–136. <https://doi.org/10.1016/j.ijpsycho.2021.04.004>
- 28 Clayson, P. E., Carbine, K. A., Baldwin, S. A., & Larson, M. J. (2019). Methodological

- 1 reporting behavior, sample sizes, and statistical power in studies of event-related  
2 potentials: Barriers to reproducibility and replicability. *Psychophysiology*, 56(11),  
3 e13437. <https://doi.org/10.1111/psyp.13437>
- 4 Clayson, P. E., Carbine, K. A., Baldwin, S. A., Olsen, J. A., & Larson, M. J. (2021). Using  
5 generalizability theory and the ERP Reliability Analysis (ERA) Toolbox for assessing  
6 test-retest reliability of ERP scores part 1: Algorithms, framework, and implementation.  
7 *International Journal of Psychophysiology: Official Journal of the International*  
8 *Organization of Psychophysiology*, 166, 174–187.  
9 <https://doi.org/10.1016/j.ijpsycho.2021.01.006>
- 10 Clayson, P. E., & Miller, G. A. (2017). Psychometric considerations in the measurement of  
11 event-related brain potentials: Guidelines for measurement and reporting. *International*  
12 *Journal of Psychophysiology: Official Journal of the International Organization of*  
13 *Psychophysiology*, 111, 57–67. <https://doi.org/10.1016/j.ijpsycho.2016.09.005>
- 14 Cowan, N. (2001). The magical number 4 in short-term memory: a reconsideration of mental  
15 storage capacity. *The Behavioral and Brain Sciences*, 24(1), 87–114; discussion 114–  
16 185. <https://doi.org/10.1017/s0140525x01003922>
- 17 Cowan, N., Elliott, E. M., Scott Saults, J., Morey, C. C., Mattox, S., Hismjatullina, A., &  
18 Conway, A. R. A. (2005). On the capacity of attention: its estimation and its role in  
19 working memory and cognitive aptitudes. *Cognitive Psychology*, 51(1), 42–100.  
20 <https://doi.org/10.1016/j.cogpsych.2004.12.001>
- 21 Cowan, N., & Morey, C. C. (2006). Visual working memory depends on attentional filtering.  
22 *Trends in Cognitive Sciences*, 10(4), 139–141.
- 23 Cowan, N. (2014). Working memory underpins cognitive development, learning, and  
24 education. *Educational Psychology Review*, 26(2), 197–223.
- 25 Daneman, M., & Carpenter, P. A. (1980). Individual differences in working memory and  
26 reading. *Journal of Verbal Learning and Verbal Behavior*, 19(4), 450–466.
- 27 de Cheveigné, A. (2020). ZapLine: A simple and effective method to remove power line  
28 artifacts. *NeuroImage*, 207, 116356. <https://doi.org/10.1016/j.neuroimage.2019.116356>

- 1 Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-  
2 trial EEG dynamics including independent component analysis. *Journal of Neuroscience*  
3 *Methods*, 134(1), 9–21. <https://doi.org/10.1016/j.jneumeth.2003.10.009>
- 4 Diamantopoulou, S., Poom, L., Klaver, P., & Talsma, D. (2011). Visual working memory  
5 capacity and stimulus categories: a behavioral and electrophysiological investigation.  
6 *Experimental Brain Research*, 209(4), 501–513. [https://doi.org/10.1007/s00221-011-](https://doi.org/10.1007/s00221-011-2536-z)  
7 2536-z
- 8 Drew, T., & Vogel, E. K. (2008). Neural measures of individual differences in selecting and  
9 tracking multiple moving objects. *The Journal of Neuroscience: The Official Journal of*  
10 *the Society for Neuroscience*, 28(16), 4183–4191.  
11 <https://doi.org/10.1523/JNEUROSCI.0556-08.2008>
- 12 Duval, S., & Tweedie, R. (2000). Trim and fill: A simple funnel-plot-based method of testing  
13 and adjusting for publication bias in meta-analysis. *Biometrics*, 56(2), 455–463.  
14 <https://doi.org/10.1111/j.0006-341x.2000.00455.x>
- 15 Emrich, S. M., Al-Aidroos, N., Pratt, J., & Ferber, S. (2009). Visual search elicits the  
16 electrophysiological marker of visual working memory. *PLoS One*, 4(11), e8042.  
17 <https://doi.org/10.1371/journal.pone.0008042>
- 18 Feldmann-Wüstefeld, T. (2021). Neural measures of working memory in a bilateral change  
19 detection task. *Psychophysiology*, 58(1), e13683. <https://doi.org/10.1111/psyp.13683>
- 20 Feldmann-Wüstefeld, T., Vogel, E. K., & Awh, E. (2018). Contralateral delay activity indexes  
21 working memory storage, not the current focus of spatial attention. *Journal of Cognitive*  
22 *Neuroscience*, 30(8), 1185–1196. [https://doi.org/10.1162/jocn\\_a\\_01271](https://doi.org/10.1162/jocn_a_01271)
- 23 Feuerstahler, L. M., Luck, S. J., MacDonald, A., 3rd, & Waller, N. G. (2019). A note on the  
24 identification of change detection task models to measure storage capacity and  
25 attention in visual working memory. *Behavior Research Methods*, 51(3), 1360–1370.  
26 <https://doi.org/10.3758/s13428-018-1082-z>
- 27 Forsberg, A., Adams, E. J., & Cowan, N. (2023). Why does visual working memory ability  
28 improve with age: More objects, more feature detail, or both? A registered report.

- 1        *Developmental Science*, 26(2), e13283. <https://doi.org/10.1111/desc.13283>
- 2        Fukuda, K., & Vogel, E. K. (2011). Individual differences in recovery time from attentional  
3        capture. *Psychological Science*, 22(3), 361–368.  
4        <https://doi.org/10.1177/0956797611398493>
- 5        Garrett-Ruffin, S., Hindash, A. C., Kaczurkin, A. N., Mears, R. P., Morales, S., Paul, K.,  
6        Pavlov, Y. G., & Keil, A. (2021). Open science in psychophysiology: An overview of  
7        challenges and emerging solutions. *International Journal of Psychophysiology: Official*  
8        *Journal of the International Organization of Psychophysiology*, 162, 69–78.  
9        <https://doi.org/10.1016/j.ijpsycho.2021.02.005>
- 10       Gathercole, S. E., & Baddeley, A. D. (1989). Evaluation of the role of phonological STM in  
11       the development of vocabulary in children: A longitudinal study. *Journal of Memory and*  
12       *Language*, 28(2), 200–213.
- 13       Gelman, A. (2007). Struggles with survey weighting and regression modeling. *Statistical*  
14       *Science: A Review Journal of the Institute of Mathematical Statistics*, 22(2), 153–164.  
15       <https://doi.org/10.1214/088342306000000691>
- 16       Gelman, A., & Loken, E. (2023). *The garden of forking paths: Why multiple comparisons can*  
17       *be a problem, even when there is no “fishing expedition” or “p-hacking” and the*  
18       *research hypothesis was posited ahead of time*. Retrieved August 21, 2023, from  
19       <http://stat.columbia.edu/~gelman/research/unpublished/forking.pdf>
- 20       Götz, M., Sarma, A., & O’Boyle, E. H. (2024). The multiverse of universes: A tutorial to plan,  
21       execute and interpret multiverses analyses using the R package multiverse.  
22       *International Journal of Psychology: Journal International de Psychologie*, 59(6), 1003–  
23       1014. <https://doi.org/10.1002/ijop.13229>
- 24       Hakim, N., Adam, K. C. S., Gunseli, E., Awh, E., & Vogel, E. K. (2019). Dissecting the neural  
25       focus of attention reveals distinct processes for spatial attention and object-based  
26       storage in visual working memory. *Psychological Science*, 30(4), 526–540.  
27       <https://doi.org/10.1177/0956797619830384>
- 28       Hakim, N., Feldmann-Wüstefeld, T., Awh, E., & Vogel, E. K. (2020). Perturbing neural

- 1 representations of working memory with task-irrelevant interruption. *Journal of Cognitive*  
2 *Neuroscience*, 32(3), 558–569. [https://doi.org/10.1162/jocn\\_a\\_01481](https://doi.org/10.1162/jocn_a_01481)
- 3 Hakim, N., Feldmann-Wüstefeld, T., Awh, E., & Vogel, E. K. (2021). Controlling the flow of  
4 distracting information in working memory. *Cerebral Cortex (New York, N.Y.: 1991)*,  
5 31(7), 3323–3337. <https://doi.org/10.1093/cercor/bhab013>
- 6 Heuer, A., & Schubö, A. (2016). The focus of attention in visual working memory: Protection  
7 of focused representations and its individual variation. *PloS One*, 11(4), e0154228.  
8 <https://doi.org/10.1371/journal.pone.0154228>
- 9 Jennings, J. R., & Wood, C. C. (1976). Letter: The epsilon-adjustment procedure for  
10 repeated-measures analyses of variance. *Psychophysiology*, 13(3), 277–278.  
11 <https://doi.org/10.1111/j.1469-8986.1976.tb00116.x>
- 12 Jongbloed-Pereboom, M., Nijhuis-van der Sanden, M. W. G., & Steenbergen, B. (2019).  
13 Explicit and implicit motor sequence learning in children and adults; the role of age and  
14 visual working memory. *Human Movement Science*, 64, 1–11.  
15 <https://doi.org/10.1016/j.humov.2018.12.007>
- 16 Kang, M.-S., & Woodman, G. F. (2014). The neurophysiological index of visual working  
17 memory maintenance is not due to load dependent eye movements. *Neuropsychologia*,  
18 56, 63–72. <https://doi.org/10.1016/j.neuropsychologia.2013.12.028>
- 19 Klaver, P., Talsma, D., Wijers, A. A., Heinze, H. J., & Mulder, G. (1999). An event-related  
20 brain potential correlate of visual short-term memory. *Neuroreport*, 10(10), 2001–2005.  
21 <https://doi.org/10.1097/00001756-199907130-00002>
- 22 Klug, M., & Kloosterman, N. A. (2022). Zapline-plus: A Zapline extension for automatic and  
23 adaptive removal of frequency-specific noise artifacts in M/EEG. *Human Brain Mapping*,  
24 43(9), 2743–2758. <https://doi.org/10.1002/hbm.25832>
- 25 Kothe, C. A., & Makeig, S. (2013). BCILAB: a platform for brain-computer interface  
26 development. *Journal of Neural Engineering*, 10(5), 056014.  
27 <https://doi.org/10.1088/1741-2560/10/5/056014>
- 28 Kundu, B., Sutterer, D. W., Emrich, S. M., & Postle, B. R. (2013). Strengthened effective

- 1 connectivity underlies transfer of working memory training to tests of short-term memory  
2 and attention. *The Journal of Neuroscience: The Official Journal of the Society for*  
3 *Neuroscience*, 33(20), 8705–8715. <https://doi.org/10.1523/JNEUROSCI.5565-12.2013>
- 4 Kuo, B.-C., Stokes, M. G., & Nobre, A. C. (2012). Attention modulates maintenance of  
5 representations in visual short-term memory. *Journal of Cognitive Neuroscience*, 24(1),  
6 51–60. [https://doi.org/10.1162/jocn\\_a\\_00087](https://doi.org/10.1162/jocn_a_00087)
- 7 Lakens, D. (2017). TOSTER: Two one-sided tests (TOST) equivalence testing. *R Package*  
8 *Version 0. 2, 5*, 648.
- 9 Lefebvre, C., Vachon, F., Grimault, S., Thibault, J., Guimond, S., Peretz, I., Zatorre, R. J., &  
10 Jolicœur, P. (2013). Distinct electrophysiological indices of maintenance in auditory and  
11 visual short-term memory. *Neuropsychologia*, 51(13), 2939–2952.  
12 <https://doi.org/10.1016/j.neuropsychologia.2013.08.003>
- 13 Leonard, C. J., Kaiser, S. T., Robinson, B. M., Kappenman, E. S., Hahn, B., Gold, J. M., &  
14 Luck, S. J. (2013). Toward the neural mechanisms of reduced working memory capacity  
15 in schizophrenia. *Cerebral Cortex (New York, N.Y.: 1991)*, 23(7), 1582–1592.  
16 <https://doi.org/10.1093/cercor/bhs148>
- 17 Liesefeld, H. R., & Müller, H. J. (2019). Current directions in visual working memory  
18 research: An introduction and emerging insights. *British Journal of Psychology (London,*  
19 *England: 1953)*, 110(2), 193–206. <https://doi.org/10.1111/bjop.12377>
- 20 Lotfi, S., Ward, R., Mathew, A., Shokoohi-Yekta, M., Rostami, R., Motamed-Yeganeh, N.,  
21 Christine, C., & Lee, H.-J. (2022). Limited visual working memory capacity in children  
22 with dyslexia: An ERP study. *NeuroRegulation*, 9(2), 98–109.  
23 <https://doi.org/10.15540/nr.9.2.98>
- 24 Luck, S. J. (2014). *An introduction to the event-related potential technique* (2nd ed.).  
25 Bradford Books.
- 26 Luck, S. J., & Gaspelin, N. (2017). How to get statistically significant effects in any ERP  
27 experiment (and why you shouldn't). *Psychophysiology*, 54(1), 146–157.  
28 <https://doi.org/10.1111/psyp.12639>

- 1 Luck, S. J., Stewart, A. X., Simmons, A. M., & Rhemtulla, M. (2021). Standardized  
2 measurement error: A universal metric of data quality for averaged event-related  
3 potentials. *Psychophysiology*, *58*(6), e13793. <https://doi.org/10.1111/psyp.13793>
- 4 Luck, S. J., & Vogel, E. K. (2013). Visual working memory capacity: from psychophysics and  
5 neurobiology to individual differences. *Trends in Cognitive Sciences*, *17*(8), 391–400.  
6 <https://doi.org/10.1016/j.tics.2013.06.006>
- 7 Luria, R., Balaban, H., Awh, E., & Vogel, E. K. (2016). The contralateral delay activity as a  
8 neural measure of visual working memory. *Neuroscience and Biobehavioral Reviews*,  
9 *62*, 100–108. <https://doi.org/10.1016/j.neubiorev.2016.01.003>
- 10 McCollough, A. W., Machizawa, M. G., & Vogel, E. K. (2007). Electrophysiological measures  
11 of maintaining representations in visual working memory. *Cortex; a Journal Devoted to*  
12 *the Study of the Nervous System and Behavior*, *43*(1), 77–94.  
13 [https://doi.org/10.1016/s0010-9452\(08\)70447-7](https://doi.org/10.1016/s0010-9452(08)70447-7)
- 14 Martin, R. C., & Romani, C. (1994). Verbal working memory and sentence comprehension: A  
15 multiple-components view. *Neuropsychology*, *8*(4), 506–523.
- 16 Miyake, A., & Shah, P. (Eds.). (1999). *Models of working memory: Mechanisms of active*  
17 *maintenance and executive control* (Vol. 506). Cambridge University Press.  
18 <https://doi.org/10.1017/cbo9781139174909>
- 19 Naveh-Benjamin, M., & Cowan, N. (2023). The roles of attention, executive function and  
20 knowledge in cognitive ageing of working memory. *Nature Reviews Psychology*, *2*(3),  
21 151–165. <https://doi.org/10.1038/s44159-023-00149-0>
- 22 Ngiam, W. X. Q., Adam, K. C. S., Quirk, C., Vogel, E. K., & Awh, E. (2021). Estimating the  
23 statistical power to detect set-size effects in contralateral delay activity.  
24 *Psychophysiology*, *58*(5), e13791. <https://doi.org/10.1111/psyp.13791>
- 25 Nunnally, J. C., & Bernstein, I. H. (1994). *Psychometric theory* (3rd ed.). McGraw-Hill.
- 26 Oldfield, R. C. (1971). The assessment and analysis of handedness: the Edinburgh  
27 inventory. *Neuropsychologia*, *9*(1), 97–113. [https://doi.org/10.1016/0028-](https://doi.org/10.1016/0028-3932(71)90067-4)  
28 [3932\(71\)90067-4](https://doi.org/10.1016/0028-3932(71)90067-4)

- 1   Olivers, C. N. L. (2008). Interactions between visual working memory and visual attention.  
2       Frontiers in Bioscience, 13(13), 1182–1191.
- 3   Pashler, H. (1988). Familiarity and visual change detection. *Perception & Psychophysics*,  
4       44(4), 369–378. <https://doi.org/10.3758/bf03210419>
- 5   Paul, M., Govaart, G. H., & Schettino, A. (2021). Making ERP research more transparent:  
6       Guidelines for preregistration. *International Journal of Psychophysiology: Official Journal*  
7       *of the International Organization of Psychophysiology*, 164, 52–63.  
8       <https://doi.org/10.1016/j.ijpsycho.2021.02.016>
- 9   Pavlov, Y. G., Adamian, N., Appelhoff, S., Arvaneh, M., Benwell, C. S. Y., Beste, C., Bland,  
10    A. R., Bradford, D. E., Bublatzky, F., Busch, N. A., Clayson, P. E., Cruse, D.,  
11    Czeszumski, A., Dreber, A., Dumas, G., Ehinger, B., Ganis, G., He, X., Hinojosa, J. A.,  
12    ... Mushtaq, F. (2021). #EEGManyLabs: Investigating the replicability of influential EEG  
13    experiments. *Cortex; a Journal Devoted to the Study of the Nervous System and*  
14    *Behavior*, 144, 213–229. <https://doi.org/10.1016/j.cortex.2021.03.013>
- 15   Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: transforming  
16    numbers into movies. *Spatial Vision*, 10(4), 437–442.  
17    <https://doi.org/10.1163/156856897x00366>
- 18   Perron, R., Lefebvre, C., Robitaille, N., Brisson, B., Gosselin, F., Arguin, M., & Jolicoeur, P.  
19    (2009). Attentional and anatomical considerations for the representation of simple  
20    stimuli in visual short-term memory: evidence from human electrophysiology.  
21    *Psychological Research*, 73(2), 222–232. <https://doi.org/10.1007/s00426-008-0214-y>
- 22   Quirk, C., Adam, K. C. S., & Vogel, E. K. (2020). No evidence for an object working memory  
23    capacity benefit with extended viewing time. *eNeuro*, 7(5), ENEURO.0150–20.2020.  
24    <https://doi.org/10.1523/ENEURO.0150-20.2020>
- 25   Roy, Y., & Faubert, J. (2023). Is the Contralateral Delay Activity (CDA) a robust neural  
26    correlate for Visual Working Memory (VWM) tasks? A reproducibility study.  
27    *Psychophysiology*, 60(2), e14180. <https://doi.org/10.1111/psyp.14180>
- 28   Sarma, A., Hwang, K., Hullman, J., & Kay, M. (2024). Milliways: Taming multiverses through

- 1        principled evaluation of data analysis paths. *Proceedings of the CHI Conference on*  
2        *Human Factors in Computing Systems*, 1–15. <https://doi.org/10.1145/3613904.3642375>
- 3        Schneider, D., Barth, A., Getzmann, S., & Wascher, E. (2017). On the neural mechanisms  
4        underlying the protective function of retroactive cuing against perceptual interference:  
5        Evidence by event-related potentials of the EEG. *Biological Psychology*, 124, 47–56.  
6        <https://doi.org/10.1016/j.biopsycho.2017.01.006>
- 7        Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations stabilize?  
8        *Journal of Research in Personality*, 47(5), 609–612.  
9        <https://doi.org/10.1016/j.jrp.2013.05.009>
- 10        Shavelson, R. J., & Webb, N. M. (2012). *Generalizability Theory*. SAGE Publications.
- 11        Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology:  
12        undisclosed flexibility in data collection and analysis allows presenting anything as  
13        significant. *Psychological Science*, 22(11), 1359–1366.  
14        <https://doi.org/10.1177/0956797611417632>
- 15        Simonsohn, U., Simmons, J. P., & Nelson, L. D. (2015). Specification curve: Descriptive and  
16        inferential statistics on all reasonable specifications. *SSRN Electronic Journal*.  
17        <https://doi.org/10.2139/ssrn.2694998>
- 18        Sterne, J. A. C., & Egger, M. (2001). Funnel plots for detecting bias in meta-analysis. *Journal*  
19        *of Clinical Epidemiology*, 54(10), 1046–1055. <https://doi.org/10.1016/S0895->  
20        4356(01)00377-8
- 21        Störmer, V. S., Li, S.-C., Heekeren, H. R., & Lindenberger, U. (2013). Normative shifts of  
22        cortical mechanisms of encoding contribute to adult age differences in visual-spatial  
23        working memory. *NeuroImage*, 73, 167–175.  
24        <https://doi.org/10.1016/j.neuroimage.2013.02.004>
- 25        Strzelczyk, D., Clayson, P. E., Sigurdardottir, H. M., Mushtaq, F., Pavlov, Y. G., Devillez, H.,  
26        Lukashevich, A., Rocha, H. A., Chung, Y. H., Ortego, K. M., Störmer, V. S., Garcia  
27        Alanis, J. C., Löffler, C., Schubert, A.-L., Biel, A. L., Tretow, A., Xu, W., Hamalainen, J.,

- 1 Lu, Z., ... et al. (2023). Contralateral delay activity as a marker of visual working  
2 memory capacity: a multi-site registered replication. In *PsyArXiv Preprints* (Issue  
3 shdea). University of Zurich.  
4 [https://scholar.google.com/citations?view\\_op=view\\_citation&hl=en&citation\\_for\\_view=x](https://scholar.google.com/citations?view_op=view_citation&hl=en&citation_for_view=x)  
5 HS491AAAAAJ:9yKSN-GCB0IC
- 6 Sundre, D. L. (1993). Book reviews : Generalizability theory: A primer, by Richard J.  
7 shavelson and Noreen M. webb. Newbury Park, CA: Sage publications, 1991,137 pp.  
8 *Evaluation Practice*, 14(2), 207–209. <https://doi.org/10.1177/109821409301400219>
- 9 Tröndle, M., & Langer, N. (2024). Decomposing neurophysiological underpinnings of age-  
10 related decline in visual working memory. *Neurobiology of Aging*, 139, 30–43.  
11 <https://doi.org/10.1016/j.neurobiolaging.2024.03.004>
- 12 Trübtschek, D., Yang, Y., Gianelli, C., Cesnaite, E., Fischer, N. L., Vinding, M. C., Marshall,  
13 T. R., Algermissen, J., Pascarella, A., Puoliväli, T., Vitale, A., Busch, N. A., & Nilsonne,  
14 G. (2023). EEGManyPipelines: a large-scale, grassroots multi-analyst study of  
15 electroencephalography analysis practices in the wild. *Journal of Cognitive*  
16 *Neuroscience*, 36(2), 217–224. [https://doi.org/10.1162/jocn\\_a\\_02087](https://doi.org/10.1162/jocn_a_02087)
- 17 Tsubomi, H., Fukuda, K., Watanabe, K., & Vogel, E. K. (2013). Neural limits to representing  
18 objects still within view. *The Journal of Neuroscience: The Official Journal of the Society*  
19 *for Neuroscience*, 33(19), 8257–8263. <https://doi.org/10.1523/JNEUROSCI.5348->  
20 12.2013
- 21 Unsworth, N., Fukuda, K., Awh, E., & Vogel, E. K. (2015). Working memory delay activity  
22 predicts individual differences in cognitive abilities. *Journal of Cognitive Neuroscience*,  
23 27(5), 853–865. [https://doi.org/10.1162/jocn\\_a\\_00765](https://doi.org/10.1162/jocn_a_00765)
- 24 Viechtbauer, W. (2010). Conducting Meta-Analyses in R with the metafor Package. *Journal*  
25 *of Statistical Software*, 36(3). <https://doi.org/10.18637/jss.v036.i03>
- 26 Villena-González, M., Rubio-Venegas, I., & López, V. (2020). Data from brain activity during  
27 visual working memory replicates the correlation between contralateral delay activity  
28 and memory capacity. *Data in Brief*, 28(105042), 105042.

- 1 <https://doi.org/10.1016/j.dib.2019.105042>
- 2 Vogel, E. K., Luck, S. J., & Shapiro, K. L. (1998). Electrophysiological evidence for a  
3 postperceptual locus of suppression during the attentional blink. *Journal of Experimental*  
4 *Psychology. Human Perception and Performance*, 24(6), 1656–1674.  
5 <https://doi.org/10.1037//0096-1523.24.6.1656>
- 6 Vogel, E. K., & Machizawa, M. G. (2004). Neural activity predicts individual differences in  
7 visual working memory capacity. *Nature*, 428(6984), 748–751.  
8 <https://doi.org/10.1038/nature02447>
- 9 von Bastian, C. C., Belleville, S., Udale, R. C., Reinhartz, A., Essounni, M., & Strobach, T.  
10 (2022). Mechanisms underlying training-induced cognitive change. *Nature Reviews*  
11 *Psychology*, 1(1), 30–41. <https://doi.org/10.1038/s44159-021-00001-3>
- 12 Wang, J., Huo, S., Wu, K. C., Mo, J., Wong, W. L., & Maurer, U. (2022). Behavioral and  
13 neurophysiological aspects of working memory impairment in children with dyslexia.  
14 *Scientific Reports*, 12(1), 12571. <https://doi.org/10.1038/s41598-022-16729-8>
- 15 Westheimer, G. (1954a). Eye movement responses to a horizontally moving visual stimulus.  
16 *A.M.A. Archives of Ophthalmology*, 52(6), 932–941.  
17 <https://doi.org/10.1001/archopht.1954.00920050938013>
- 18 Westheimer, G. (1954b). Mechanism of saccadic eye movements. *A.M.A. Archives of*  
19 *Ophthalmology*, 52(5), 710–724.  
20 <https://doi.org/10.1001/archopht.1954.00920050716006>
- 21 Widmann, A., & Schröger, E. (2012). Filter effects and filter artifacts in the analysis of  
22 electrophysiological data. *Frontiers in Psychology*, 3, 233.  
23 <https://doi.org/10.3389/fpsyg.2012.00233>
- 24 Wilkinson, G. N., & Rogers, C. E. (1973). Symbolic description of factorial models for  
25 analysis of variance. *Journal of the Royal Statistical Society. Series C, Applied*  
26 *Statistics*, 22(3), 392.
- 27 Williams, J. R., Robinson, M. M., Schurgin, M. W., Wixted, J. T., & Brady, T. F. (2022). You  
28 cannot “count” how many items people remember in visual working memory: The

1 importance of signal detection-based measures for understanding change detection  
2 performance. *Journal of Experimental Psychology. Human Perception and*  
3 *Performance*, 48(12), 1390–1409. <https://doi.org/10.1037/xhp0001055>

4

5

Journal Pre-proof

## Scientific transparency statement

DATA: All raw and processed data supporting this research are publicly available:  
[https://gin.g-node.org/EEGManyLabs/EEGManyLabs\\_Replication\\_VogelMachizawa2004\\_Raw](https://gin.g-node.org/EEGManyLabs/EEGManyLabs_Replication_VogelMachizawa2004_Raw),  
[https://gin.g-node.org/EEGManyLabs/EEGManyLabs\\_Replication\\_VogelMachizawa2004\\_Processed](https://gin.g-node.org/EEGManyLabs/EEGManyLabs_Replication_VogelMachizawa2004_Processed)

CODE: All analysis code supporting this research is publicly available:  
<https://github.com/ksgfan/EEGManyLabs>

MATERIALS: All study materials supporting this research are publicly available:  
<https://github.com/ksgfan/EEGManyLabs>

DESIGN: This article reports, for all studies, how the author(s) determined all sample sizes, all data exclusions, all data inclusion and exclusion criteria, and whether inclusion and exclusion criteria were established prior to data analysis.

PRE-REGISTRATION: At least part of the study procedures was pre-registered in a time-stamped, institutional registry prior to the research being conducted:  
<https://doi.org/10.31234/osf.io/shdea> At least part of the analysis plans was pre-registered in a time-stamped, institutional registry prior to the research being conducted:  
<https://doi.org/10.31234/osf.io/shdea> The analyses that were undertaken deviated from the preregistered analysis plans. All such deviations are fully disclosed in the manuscript.

For full details, see the *Scientific Transparency Report* in the supplementary data to the online version of this article.