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

Dawid Strzelczyk[1,2], Peter E. Clayson[3], Heida Maria Sigurdardottir[4], Faisal Mushtaq[5], Yuri G. Pavlov[6], Hélène Devillez[4], Anton Lukashevich[4], Harold A. Rocha[3], Yong Hoon Chung[7], Kevin M. Ortego[7], Viola S. Stoermer[7], José C. García Alanis[8], Christoph Löffler[8], Anna-Lena Schubert[8], Anna Lena Biel[9], Ariane Tretow[10], Weiyong Xu[10], Jarmo Hämäläinen[10], Zitong Lu[11], Yong Min Choi[11], Eva Lout[11], Julie D. Golomb[11], Shuangke Jiang[12], Myles Jones[12], Eda Mizrak[12,13], Claudia C. von Bastian[12], Niko A. Busch[9], Maro G. Machizawa[14,15], William X. Q. Ngiam[16,17], Edward K. Vogel[16,17], Nicolas Langer[1,2,18*]

[1] Methods of Plasticity Research, Department of Psychology, University of Zurich, Zurich, Switzerland
[2] Neuroscience Center Zurich (ZNZ), University of Zurich and ETH Zurich, Zurich, Switzerland
[3] University of South Florida, Tampa, Florida, USA
[4] University of Iceland, Reykjavik, Iceland
[5] University of Leeds, Leeds, UK
[6] University of Tuebingen, Tuebingen, Germany
[7] Dartmouth College, Hanover, New Hampshire, USA
[8] University of Mainz, Germany
[9] University of Münster, Germany
[10] University of Jyväskylä, Finland
[11] The Ohio State University, Columbus, Ohio, USA
[12] University of Sheffield, Sheffield, UK
[13] University of Oxford, Oxford, UK
[14] Xiberlinc Inc., Tokyo, Japan
[15] Hiroshima University, Hiroshima, Japan
[16] Department of Psychology, University of Chicago, Illinois, USA
[17] Institute of Mind and Biology, University of Chicago, Illinois, USA
[18] Center of Reproducible Science, University of Zurich, Zurich, Switzerland


* n.langer@psychologie.uzh.ch

# Abstract

Visual working memory (VWM) is a temporary storage system capable of retaining information that can be accessed and manipulated by higher cognitive processes, thereby facilitating a wide range of cognitive functions. Electroencephalography (EEG) is used to understand the neural correlates of VWM with high temporal precision, and one commonly used EEG measure is an event-related potential called the contralateral delay activity (CDA). In a landmark study by Vogel and Machizawa (2004), the authors found that the CDA amplitude increases with the number of items stored in VWM and plateaus around three to four items, which is thought to represent the typical adult working memory capacity. Critically, this study also showed that the increase in CDA amplitude between two-item and four-item arrays correlated with individual subjects' VWM performance. Although these results have been supported by subsequent studies, a recent study suggested that the number of subjects used in experiments investigating the CDA may not be sufficient to detect differences in set size and to provide a reliable account of the relationship between behaviorally measured VWM capacity and the CDA amplitude. To address this, the current study, as part of the #EEGManyLabs project, aims to conduct a multi-site replication of Vogel and Machizawa's (2004) seminal study on a large sample of participants, with a pre-registered analysis plan. Through this, our goal is to contribute to deepening our understanding of the neural correlates of visual working memory.

# Introduction

Visual working memory (VWM) is a temporary storage system that holds information that can be accessed and manipulated by higher cognitive functions (Luck & Vogel, 2013). Visual working memory is considered a central construct in cognitive neuroscience and is a putative intermediary for information transfer (Liesefeld & Müller, 2019), thereby facilitating various cognitive functions, including reading comprehension (Lotfi et al., 2022; Wang et al., 2022), planning and problem-solving (Cowan et al., 2005; Miyake & Shah, 1999; Naveh-Benjamin & Cowan, 2023), and learning new skills (Jongbloed-Pereboom et al., 2019; Pickering, 2006; von Bastian et al., 2022).

There is a large body of work employing electroencephalography (EEG) to understand the neural correlates of VWM with high temporal precision. A commonly used EEG measure of VWM is an event-related potential (ERP) called the contralateral delay activity (CDA). This signal has also been referred to in other studies as Contralateral Negative Slow Wave (CNSW; Klaver et al., 1999), Sustained Posterior Contralateral Negativity (SPCN; Brisson & Jolicoeur, 2007 and Perron et al., 2009), and Contralateral Search Activity (CSA; Emrich et al., 2009). These different terms all refer to the same visual working memory correlate. Hence, we will maintain the use of the term CDA throughout the remainder of the paper.

The CDA is a difference wave constructed by contrasting activity ipsilateral and contralateral to to-be-remembered items. In a typical experiment on the CDA, items are shown bilaterally but only those on one side of the screen are supposed to be memorized. The idea of the subtraction is to eliminate any activity related to early perceptual and low-level processing by assuming that they equally affect ipsi- and contralateral ERPs. Activity over the contralateral hemisphere tends to be more negative than ipsilateral activity during VWM retention (Luria et al., 2016; Ngiam et al., 2021; Vogel & Machizawa, 2004). Thus, it has been suggested that the CDA reflects the neural activity related to maintenance of information in visual working memory, and studies have shown that the amplitude and duration of the CDA are linked to the amount of information stored in working memory.

In a seminal paper from Vogel and Machizawa (2004), the authors demonstrated that CDA amplitude increases with the number of items stored in VWM and plateaus at around three to four items, consistent with the typical adult working memory capacity (Cowan, 2001; Forsberg et al., 2023; Pashler, 1988). More importantly, Vogel and Machizawa (2004) showed that the increase in CDA amplitude with greater memory load correlated with individual VWM performance. Specifically, individuals with high VWM capacity exhibited a larger increase in CDA when attempting to memorize four compared to two items than individuals with low VWM capacity. In this study, the CDA was elicited using a color change detection task. The task involves presenting participants with a central arrow cue that indicates whether participants need to memorize items on the left or right of screen center. The cue is followed by a bilateral stimulus array with equal numbers of colored squares shown on each side (set size 1 to 10). After a short retention phase, participants are presented with a second array and asked to indicate whether any of the squares on the cued side changed color (Figure 2). The lateralized color change detection task is now a widely used paradigm to examine visual working memory processes (for a review see: Luria et al., 2016; Feldmann-Wüstefeld, 2021) and has been explored in several variations such as different set sizes, including distractions, retro-cueing and using different shapes and colors (Feldmann-Wüstefeld, 2021; Feuerstahler et al., 2019; Luria et al., 2016; Roy & Faubert, 2023; Schneider et al., 2017).

The finding that the CDA amplitude is sensitive to the number of stimuli to be remembered has been replicated in numerous studies (Asp et al., 2021; Brady et al., 2016; Hakim et al., 2019; Heuer & Schubö, 2016; Quirk et al., 2020; Unsworth et al., 2015). Furthermore, several studies have validated the positive correlation between CDA amplitude increase and VWM capacity (Adam et al., 2018; Feldmann-Wüstefeld et al., 2018; Villena-González et al., 2020). In the review paper of Luria et al. (2016), the authors conducted a meta-analysis from 11 previous studies and reported an aggregated correlation of $r = .596$ (Luria et al., 2016). However, a recent study indicated that the typical numbers of subjects and trials for CDA experiments seen in the literature may be underpowered for detecting set size differences (Ngiam et al., 2021).

The insufficient power issue is even more pressing for the correlation between the VWM capacity and the CDA amplitude increase. Critically, Schönbrodt and Perugini (2013) demonstrated that correlation estimates typically stabilize at a sample size of approximately 250 subjects. Except for one large study (N = 171; Unsworth et al., 2015), the average sample size of previous studies investigating the relationship between VWM capacity and CDA amplitude was 32 subjects (range 12-83 subjects for 12 studies; Luria et al., 2016). Finally, the inherent flexibility in EEG analysis, including analysis of the CDA, leaves many decisions up to the researcher. This leaves open the possibility to exploit these "researcher's degrees of freedom" (i.e., garden of forking paths; Gelman & Loken, 2013), either intentionally or unintentionally. Such practices to find statistically significant effects can lead to erroneous inferences and perpetuate replication problems in cognitive neuroscience (Clayson et al., 2019; Luck & Gaspelin, 2017).

To address these general issues in the EEG literature, a new initiative, #EEGManyLabs (Pavlov et al., 2021), was recently launched. #EEGManyLabs is an international, collaborative effort focussed on directly replicating some of the most influential EEG studies published to establish the robust of, and confidence in, widely cited phenomena. Importantly, the #EEGManyLabs project is designed to address some of the limitations of previous replication efforts by using a large sample of participants, standardized procedures, and a pre-registered analysis plan (i.e., Registered Report; Pavlov et al., 2021).

Given the aforementioned relatively recent meta-analysis on CDA in VWM (Luria et al., 2016), it is worth noting the specific advantages conferred by large-scale direct replication over meta-analyses. While a meta-analysis provides a valuable summary of published, or publicly available, research, it is highly vulnerable to the issue of publication bias (Levine et al., 2009; Lin & Chu, 2018; Rothstein et al., 2005; Yang et al., 2023), threatening its validity and a risk of over-estimation of effect sizes and increase in the rate of false positives (Bartoš et al., 2023; Lane & Dunlap, 1978; Yang et al., 2023). In contrast, a large-scale direct replication, as operationalised in #EEGManyLabs, has the advantage of undertaking a new study on a larger

sample, increasing statistical power, and enhancing the robustness of the results. The adoption of identical materials across all participating laboratories in a large-scale replication ensures consistency and constrains potentially confounding sources of variation and measurement error. Consequently, any observed heterogeneity in the results is primarily attributed to differences in the subject populations and site-specific differences across the laboratories, allowing for a clearer understanding of the true effects being examined.

As part of the #EEGManyLabs project, the current study aims to contribute to the existing literature on VWM and the CDA by conducting a robust multi-site, large-scale replication of Vogel and Machizawa's (2004) seminal study. The present study was chosen for replication by a global consortium of EEG specialists owing to its scientific significance (for further information on the selection process, refer to Pavlov et al., 2021). In accordance with the #EEGManyLabs project, this Registered Report will closely adhere to the original study design and ensure adequate statistical power with a large sample size. The present study will also follow preregistered analysis steps to ensure the integrity of the direct replication and statistical inferences (Paul et al., 2021).

In line with the original study, the following hypotheses will be tested:

The investigation of the relationship between the CDA amplitude and the number of items stored in memory will be conducted by replicating Experiment #3 of the original study, using three set sizes (2, 4, and 6 items per side).

[H1.1] CDA amplitude increases from arrays of two items per side to arrays of four items per side.

[H1.2] CDA amplitude increases from arrays of two items per side to arrays of six items per side.

[H1.3] CDA amplitude for four items and six items are equivalent.

Additionally, the study will examine whether the CDA amplitude is related to performance on the change detection task:

[H2.1] VWM capacity (measured behaviorally) is positively correlated with CDA amplitude increase from two to four items.

[H2.2] Subject's VWM capacity (measured behaviorally) is not correlated with CDA amplitude increase from four to six items.

Finally, replication success is established for each hypothesis separately. It is operationally defined as a statistically significant random-effects meta-analytic estimate in the same direction as in the original study or as a null effect, depending on the predictions of the respective hypotheses. These outcomes are obtained by combining results from all laboratories.

# Method

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 1 Registered Report that follows guidelines for open science in psychophysiological research as outlined by (Garrett-Ruffin et al., 2021). The study materials and initial code for data processing are available on the Open Science Framework (OSF) at the link https://osf.io/pbr8c/. Full information about each site (i.e., EEG, recruitment) will be also posted on OSF. The Open Neuro (https://openneuro.org/) repository will provide access to raw EEG data. Each site has obtained approval from the local ethics committee to conduct the study and share data.

## Known differences from the original study:

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

| Offline Processing Step | Vogel & Machizawa (2004) | #EEGManyLabs direct replication pipeline | #EEGManyLabs advanced pre-processing pipeline |
|---|---|---|---|
| Filter | 1. Hardware online filter: bandpass of 0.01-80 Hz (half-power cutoff, Butterworth filters)<br>2. 35 Hz LP only for plots | 1. Offline bandpass of 0.01-80 Hz (half-power cutoff, Butterworth filters)<br>2. 35 Hz LP for plotting | 1. Offline bandpass of 0.01-80 Hz (half-power cutoff, Butterworth filters)<br>2. 35 Hz LP for plotting |
| Line noise removal | no | ZapLine method[1] | ZapLine method |
| Ocular artifact rejection | 1. Trials containing ocular artifacts were removed (i.e., blinks or eye 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).<br>2. Blinks: unipolar VEOG >50 microvolt | 1. Trials containing ocular artifacts will be removed (i.e., blinks or eye movements larger than 1 visual degree). A calibration paradigm will be used to estimate the subject specific amplitude representing 1 visual degree<br>2. Blinks: unipolar VEOG >50 microvolt | If eye-tracker recording is available, we will exclude trials with eye movements larger than 1 degree. If no eye-tracker is available, we identify ocular artifacts using a tailored subject-specific amplitude threshold for the EOG electrodes, which is obtained from the saccadic calibration task. |
| Artifact Rejection | 1. peak-to-peak amplitude >200 microvolt<br>2. Visual inspection was used to identify and exclude trials containing movement artifacts or blocking | 1. peak-to-peak amplitude >200 microvolt<br>2. Visual inspection will be used to identify and exclude trials containing movement artifacts or blocking | 1. Peak-to-peak amplitude >200 microvolt<br>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<br>2. Visual inspection | 1. peak-to-peak amplitude >75 microvolt. Bad channels will not be interpolated, but artifactual trials will be rejected<br>2. Visual inspection | 1. Correlation below 0.85 with neighbouring channels<br>2. 4 SD or more line noise relative to signal than all other channels<br>3. Blocking longer than 5 s |
| Bad Channel Interpolation | N/A | N/A | Spherical spline interpolation |
| Reference | Algebraic average of the left and right mastoids | Algebraic average of the left and right mastoids | Algebraic average of the left and right mastoids |
| Baseline Interval | -200 - 0 ms | -200 - 0 ms | -200 - 0 ms |
| Region of Interest | left electrode cluster: P3, T5/P7, O1;<br>right electrode cluster: P4, T6/P8, O2 | left electrode cluster: P3, T5/P7, O1;<br>right electrode cluster: P4, T6/P8, O2 | left electrode cluster: P3, T5/P7, O1;<br>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) | retention phase (i.e., 300 - 900 ms after the onset of memory array) |
| Set Size | Experiment #3: 2,4,6 | 2, 4, 6 | 2, 4, 6 |
| Visual Memory Capacity | K | K & d' | K & d' |

[1] Note that several labs are recording the task with an eye-tracker, which will induce line noise. Therefore, we decided to use ZapLine to reduce the line noise.

## Sample size & Inclusion criteria

Participants will be recruited from universities or nearby communities. The study will only include individuals between 18 and 35 years free from any diagnosed psychiatric or neurological disorders and with intact color vision. We will acquire demographics (i.e., age, gender), handedness (Edinburgh Handedness Inventory [EHI]; Oldfield, 1971) and education level based on International Standard Classification of Education (ISCED) (http://uis.unesco.org/en/topic/international-standard-classification-education-isced).

The required sample size is estimated for each hypothesis: For hypothesis #1, we used the CDA power calculator (https://williamngiam.shinyapps.io/CDAPower/; Ngiam et al., 2021) to estimate the required sample size to detect a set-size effect between set size 2 and 4 (which is similar as between set size 2 and 6). With a minimum of 170 clean trials per condition (i.e.,

excluding subjects with a bad trial rate > 30%) and 90% power, the estimated number of subjects required is 70 (see Figure 1).

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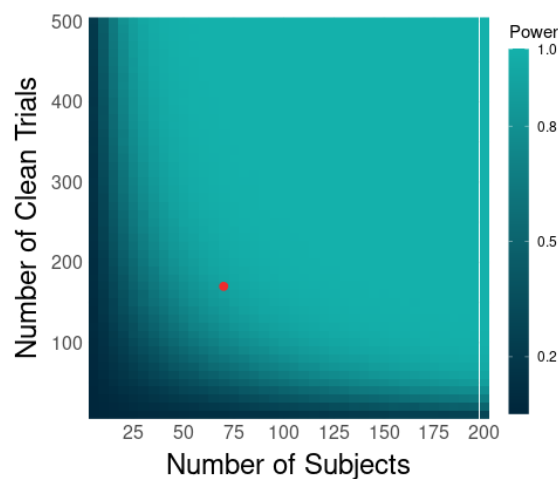


*Figure 1. CDA power calculation. We estimated the required sample size for the set size effect between set size 2 and 4, assuming at least 170 trials (i.e., maximum of 30% bad trials) and 90% power. The estimated number of subjects required is 70 (red dot).*

## Exclusion criteria

The color change detection task requires the participants to discriminate between colored squares, therefore color blindness is a critical exclusion criterion. We will test the color-vision

with an online color-vision test (https://colormax.org/color-blind-test/). If a participant has a color vision deficiency, the participant will be excluded.

**Exclusion criteria for direct replication**

Following the original study, we will exclude trials with eye movements, blinks, and blocking (amplifier saturation after drift). To identify eye movements and blinks, horizontal electrooculography (EOG) will be concurrently recorded. Contaminated trials will be identified by large (>1°) eye movements (Vogel & Machizawa, 2004).

In the original study, the authors used a heuristic for 1° horizontal eye movement and a fixed amplitude threshold for each subject. In this replication study, we will deviate from this procedure and calculate the 1° 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° horizontal eye movement, there will be a separate horizontal EOG saccade calibration task prior to the main experiment. This task involves participants making saccades to left and right targets on the screen. Participants will start each trial by fixating at the center of the screen. Following a key press there will be a jittered interval between 1200~1600 ms and a saccade target (a red disk; 0.6° in size) will appear either 3° or 6° away from the fixation on the left or right side of the screen along the horizontal midline. Participants are instructed to make a saccade to the target location as soon as it appears and press a space bar once they have successfully made the saccade. There will be 15 trials per condition, resulting in 60 trials total. The data from the saccade calibration paradigm will be 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 correction using a pre-stimulus baseline interval of -200 to 0 ms. Given previous research showing saccade onset latency being ~200 ms (Westheimer, 1954a, 1954b), horizontal EOG (i.e., HEOG = HEOGR - HEOGL) channel amplitudes from horizontal saccades will be averaged during the 300~400 ms interval across the left and right conditions. The 1° horizontal eye movement amplitude threshold will then be calculated by extrapolating from 3° and 6° eye

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

**Alternative analysis pipelines**

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

## Procedure

Upon their arrival, participants will receive a brief overview of the experiment and will be asked to give their informed written consent for participating in the study and allowing data sharing. Next, the participants will be asked to fill out a short questionnaire regarding history of

psychiatric and neurological disorders, handedness and educational level and carry out an online color-blind test (https://colormax.org/color-blind-test/).

Next, participants will be comfortably seated in a chair. If available, the experiment will be conducted in a sound- and electrically-shielded Faraday recording cage (see Table 2). The cage is equipped with a chinrest to minimize head movements. A cap with integrated electrodes will be placed on the participant's head and impedances will be checked if provided by the EEG amplifier system and improved if necessary (see Table 2 for details). As this project is part of a wider initiative on replicability in EEG (#EEGManyLabs), several of the laboratories in this replication will also collect resting state EEG recordings together with some personality measures (https://osf.io/sp3ck/; Pavlov et al., 2021). Neither resting EEG nor personality data will be analyzed in the current study but will be merged across sites as part of a future replication project to be reported elsewhere. Subsequently, the actual color change detection task experiment will begin. The expected duration of the entire experiment is approximately 120 minutes. Upon completion of the examination, participants will receive compensation or credit for their participation.

## Experimental Paradigm

The color change detection task is identical to the task used in the original study (Vogel & Machizawa, 2004). The paradigm was implemented in MATLAB, using the PsychToolbox extensions (Brainard, 1997; Pelli, 1997). Each trial of the task starts with a fixation cross presented for a random time between 300 and 400 ms. Then, a central arrow appears for 200 ms, indicating which side of the screen the participant should pay attention to. This is followed by another fixation period of a random time interval between 300 ms and 400 ms. Afterwards, a memory array will be presented for 100 ms, which will consist of either 2, 4, or 6 colored squares on each side of the screen. Participants will be instructed to only memorize the part of the memory array indicated by the arrow. This will be followed by a 900 ms retention interval with a blank screen and a fixation cross (see Figure 2). Finally, a test array will be presented for 2000 ms, and the participants will have to indicate whether the test array is identical to the

previous memory array ("no-change" trial) or whether the test array was different by one color ("change" trial). The participants will indicate whether a change occurred by pressing either the A or L button on a keyboard. The button they press will depend on the instruction they were given, with half of the participants being instructed to press the A button for a change and the L button for no change, while the other half will be instructed to do the opposite. Participants will be instructed to use their left hand to press the A button and their right hand to press the L button. During the task, participants will be asked to focus their gaze on the fixation cross in the center of the screen until the probe appears.

All stimuli will be displayed within two regions that are 4° x 7.3° in size and are located 3° to the left and right of a central fixation cross on a gray background (8.2 cd m$^{-2}$). Each memory array consists of 2, 4 or 6 colored squares (0.65° x 0.65°) in each visual field. The squares are chosen at random from a set of seven highly distinct colors (red, blue, violet, green, yellow, black and white), and a specific color appears no more than twice in a single array. In simpler terms, a color may appear in both hemifields, but it will never be repeated within a single hemifield. The positions of the stimuli are randomized on each trial, with the restriction that the distance between squares within a visual field is at least 2° (center to center). In 50% of the trials, the color of one square in the test array on the cued side is different from the corresponding square in the memory array, while in the remaining trials, the colors of the two arrays are identical.

The task is divided into five blocks, each containing 144 trials (i.e., 720 trials per subject and 240 trials per condition and subject). The cue direction (left or right) and set size (2, 4 or 6 items on each side of the screen) will be randomly varied throughout the trials to ensure a balanced distribution of all conditions in each block. In line with the original study, no training exercise will be conducted prior to the main task.
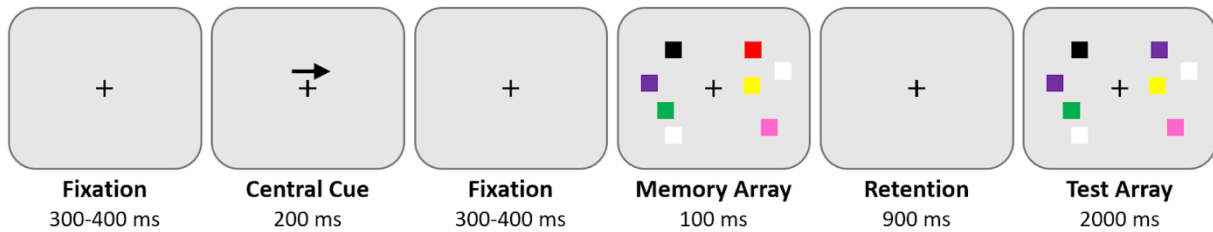
**Fixation** 300-400 ms | **Central Cue** 200 ms | **Fixation** 300-400 ms | **Memory Array** 100 ms | **Retention** 900 ms | **Test Array** 2000 ms

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

**Neurophysiological Data Acquisition**

The replicating labs will be using one of the following EEG systems and eye trackers (if available to the participating lab). Details about the acquisition setups are described in Table 2. All labs will provide the raw data to Zurich's Lab, where it will be pre-processed and analyzed. In general, large-scale studies of adults indicate that data from across sites can be combined when basic data curation and processing steps are aligned (Bigdely-Shamlo et al., 2020). In case we encounter data alignment problems, we will conduct re-centering and re-scaling techniques as introduced in previous research (Bleuzé et al., 2021; Maman et al., 2019; Mellot et al., 2023; Rodrigues et al., 2019).

*Table 2. Data acquisition settings at each lab*

| 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 |
|---|---|---|---|---|---|---|---|---|---|---|
| University of Zurich | 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 |
| 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 Sheffield | Iiyama G-master GB2488HSU; | Psychtoolbox 3.0.18 | 50cm (nasion to screen | Biosemi; 64 channels; Record at 2048 Hz | Cz; Yes | Not Available (Only offset + 25 mV) | No | Yes | Yes | No |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 531.4 x 298.9mm; 1920 x 1080; 144 Hz | | distance; | and then down-sample | | | | | | |
| 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 |
| Icelandic Vision Lab, University of Iceland | 2560*1440 60 Hz ASUS PG278QR 27" | Psychtoolbox 3.0.18 | 57cm (nasion to screen distance; | BrainVision; 32 channels; 1000 Hz | Cz, Fpz | Kept below 15 KOhm | No | Yes | No | Yes |
| The Ohio State University | BENQ XL2420-B; 1920 x 1080; 120 Hz | Psychtoolbox 3.0.18 | 80cm | BrainVision ;32 channels; 1000Hz | Cz, Fpz | Kept below 20 KOhm | EyeLink 1000; 500 Hz | Yes | Yes | Yes |
| University of Münster | ViewpixxEEG 1920 x 1080 120Hz | Psychtoolbox 3.0.18 | 86 | Biosemi; 64 + 3 channels; 1024Hz | Reference free; GND adjacent to POz | Not available with Biosemi | Eye Link, 500 Hz | Yes | no | Yes |
| University of Jyväskylä | Asus, 1920 x 1080, 120 Hz | Psychtoolbox 3.0.16 | 67 cm | Neurone; Easycap; 64 channels; 500 Hz | FCz; AFz; | Kept below 20 KOhm | EyeLink 1000; 500Hz | Yes | No | Yes |
| University of Mainz | Eizo ColorEdge CS2420; 24.1'' diag; 1920x1200; 60 Hz | Psychtoolbox 3.0.18 | 67 cm | BrainProducts; 64 channels; 1000 Hz | Cz; Fpz | Kept below 10 KOhm | No | Yes | Yes | Yes |
| North Dakota State University | ASUS ROG Strix XG27AQ 27"; 2560 x 1440; | Psychtoolbox 3.0.18 | 50 cm | Biosemi; 64 + 8 channels; 512Hz | Reference free; GND adjacent to POz | Not available with Biosemi | Eye Link 1000, 500 Hz | Yes | Yes | Yes |

## Artifact Removal and Data Pre-processing

All EEG data will be imported into EEGLAB (Delorme & Makeig, 2004) and processed using two pipelines: a pipeline that follows the original study as closely as possible (see Vogel & Machizawa, 2004), and a pipeline that utilizes a more recent analysis technique.

### Direct replication pre-processing pipeline

Following the Vogel & Machizawa (2004) study, we will downsample the data to 250 Hz and apply a bandpass filter of 0.01-80 Hz (half-power cutoff, Butterworth filters) using the EEGLAB function pop_eegfiltnew (Widmann & Schröger, 2012). Since certain labs may be measuring

eye movements using an eye tracker or may not have access to a faraday cage, line noise (50 Hz in Europe, 60 Hz in the USA) may be introduced as a result. To mitigate this, we will use ZapLine (de Cheveigné, 2020) to remove line noise artifacts by removing seven power line components. The Algorithm is highly effective at removing power line artifacts while preserving non-artifactual parts of the signal (de Cheveigné, 2020). This deviation from the original study is necessary to ensure accurate measurements. Afterwards, we will re-reference the data to an algebraic average of the left and right mastoids. We will segment the data from -200 to 1200 ms after the presentation of the memory array. The segments with saccadic eye movements (greater than 1° from the fixation cross) will be excluded from further analysis using horizontal EOG channel response data from the saccade calibration task (for detailed information please refer to *Exclusion Criteria*). Furthermore, blinks will be detected by using an amplitude threshold (>50 microvolt) in the unipolar VEOG channel. In addition, a segment will be marked as bad if any electrodes of interest (i.e., HEOG, P3, T5/P7, O1; P4, T6/P8, O2) had a peak-to-peak amplitude >75 microvolt within one time window (i.e., bad channel criteria). Finally, visual inspection will be used to identify bad trials. Finally, a baseline correction was applied using a pre-stimulus interval of -200 to 0 ms.

**Advanced pre-processing pipeline**

In addition to following the original study's data pre-processing protocol, the data will also be processed using recent advancements in neuroscience to assess the robustness of the results. First, error-prone channels will be detected by the algorithms implemented in the EEGLAB plugin clean_rawdata (http://sccn.ucsd.edu/wiki/Plugin_list_process) without applying ASR (automated subspace removal). An electrode is defined as an error-prone when recorded data from that electrode is correlated at less than 0.85 to an estimate based on neighboring electrodes. Furthermore, an electrode is defined as error-prone if it has more line noise (i.e., 50 Hz in Europe, 60 Hz in USA) relative to its signal than all other electrodes (4 standard deviations). Finally, if an electrode has a longer flat line than 5 s, it is considered error prone. These error-prone electrodes will automatically be removed and later be interpolated using a spherical spline interpolation (EEGLAB function eeg_interp.m). Next, data will be filtered using

a bandpass filter of 0.01-80 Hz (half-power cutoff, Butterworth filters). Again, we will use ZapLine (seven components) to remove line noise. Subsequently, the data will be re-referenced to an algebraic average of the left and right mastoids, segmented from -200 to 1200 ms after presentation of memory array. To determine whether a trial is artefactual, three criteria will be applied: First, we will exclude trials with blinks and large saccadic eye movements. If eye-tracker recording is available, we will exclude trials with eye movements larger than 1°. If no eye-tracker is available, we identify ocular artifacts using a tailored subject-specific amplitude threshold for the HEOG electrodes, which is obtained from the saccadic calibration task. Second, a sliding time window approach will be adopted from Adam et al. (2018). To identify trials containing blocking artifacts, a sliding time window of 200ms will be shifted across the segments without overlap. If any time window contained 30ms of flat line activity in any channel (i.e., range of amplitudes <1 microvolt), the corresponding segment was marked as bad. Third, to identify trials containing large amplitude artifacts, non-overlapping sliding time windows of 14ms will be used. A segment will be marked as bad if any electrode of interest (i.e., HEOG, P3, T5/P7, O1; P4, T6/P8, O2) had a peak-to-peak amplitude >200 microvolt within one time window. To foster scientific transparency and enable exact methodological replications and reproducibility, no visual inspection for bad trials rejection will be conducted, because this decision is subjective. Finally, a baseline correction will be applied using a pre-stimulus baseline interval of -200 to 0 ms.

## Confirmatory analysis plan

### CDA extraction

To remove contribution of any non-VWM-specific, bilateral activity, the CDA is computed as a difference wave on a trial-by-trial basis by subtracting activity ipsilateral to cued items (presented left or right of screen center) from contralateral activity. The CDA amplitude is extracted from a time window of 300–900 ms after the onset of the memory array. We will compute the mean CDA amplitude for each participant separately for each set size (i.e., 2, 4 and 6 cued items). For the computation of CDA, we use posterior parietal, lateral occipital, and

posterior temporal electrode sites (i.e., left electrode cluster: P3, T5/P7, O1; right electrode cluster: P4, T6/P8, O2). First the difference is calculated in electrode pairs (P3/P4), (P7/P8), (O1/O2) and then averaged. The CDA is calculated on a trial-by-trial basis for all set sizes and for all trials (i.e., correct and incorrect trials). The final step is to compute the overall average CDA for each set size by averaging the CDA of the right and left cue direction of the respective set size.

**Data Quality and Psychometric Internal Consistency**

Estimates of data quality and psychometric internal consistency will be reported. Data quality characterizes measurement error, and psychometric internal consistency provides information about whether measurement error is low enough to discriminate scores between people, which is crucial for studying individual differences (Clayson, Brush, et al., 2021; Clayson & Miller, 2017b; Luck et al., 2021). These metrics are reported to characterize the obtained data, but data will not be excluded based on these metrics to be consistent with the procedures of the original study.

Arithmetically derived estimates of the standard error of the mean will be used to characterize data quality (Luck et al., 2021). These estimates will separately quantify the precision of CDA for each set size (2, 4, and 6 cued items) using single-trial estimates of CDA (contralateral-ipsilateral activity differences).

Psychometric internal consistency estimates will use generalizability theory equations to compute coefficients of dependability for difference scores (Baldwin et al., 2015; Brennan, 2005; Clayson, Baldwin, et al., 2021; Clayson, Brush, et al., 2021; Sundre, 1993). Time-window mean amplitude estimates of single-trial trial scores of ipsilateral and contralateral activities will be used to estimate the observed group-level internal consistency of the difference scores. Dependability of contralateral-ipsilateral activity difference scores will be estimated separately for each set size and data collection site using the ERP Reliability Analysis Toolbox (Clayson, Carbine, et al., 2021; Clayson & Miller, 2017a). Because CDA scores are calculated as the difference between activity from different electrode sites on the same trial, residual covariances

will be estimated because the constituent events of the difference scores are co-occurring (Clayson, Baldwin, et al., 2021).

## Outcome-neutral test

To ensure that the data can test the stated hypotheses, we are including quality checks (see also pilot section) and outcome-neutral tests. As outcome-neutral test, we test the presence of an asymmetry between contra- or ipsilateral electrode clusters time-locked to the memory array. For this, we will average the event-related potentials across all set sizes and all subjects (i.e., grand averaged ERP) elicited by memory arrays that were either contra- or ipsilateral to electrode position. A paired sample *t* test for CDA between ipsilateral and contralateral sites will be performed separately at each study site to verify the expected within-lab CDA experimental effect. If the *t* test is significant ($p < .05$) with more negative CDA for contralateral activity than for ipsilateral activity, then this pattern of effect will justify moving forward with testing the proposed hypotheses.

### Statistical analysis

For all the statistical analyses, frequentist and Bayesian approaches will be employed. To estimate effect sizes, the statistical analysis will be initially conducted for each participating lab separately. Because of the small sample size in each lab, we will refrain from interpretation of the lab specific statistics. However, the overall replication success for the project will be determined based on meta-analytic pooled effect sizes, as per the defined criteria.

### Statistical analysis for Hypothesis #1

A repeated-measures ANOVA of the CDA amplitude in the original study revealed a significant main effect for 'set size'. Post-hoc t-tests showed significant increases in CDA amplitude for set sizes 4 and 6 compared to set size 2, with no significant differences between set sizes 4 and 6. In accordance with the original study, we will also conduct repeated-measures ANOVA. The significance level will be set to $p < 0.02$ uncorrected for multiple comparisons. If the ANOVA reveals a significant main effect, we will further conduct post-hoc t-tests (with a

significance level of p < 0.02, one-sided). We specifically test one-sided, because we hypothesize a significant increase in CDA amplitude from arrays of two items per side to arrays of four items per side [H1.1] and six items per side [H1.2.]. As in the original study, we will adjust the p-values with the Greenhouse-Geisser correction for nonsphericity (Jennings & Wood, 1976). If the ANOVA reveals a significant main effect, and the post-hoc t-tests show a significant increase in CDA amplitude between arrays of two items per side and arrays of four items per side or six items per side, it will support hypotheses [H.1.1] and [H.1.2], respectively. We will run the corresponding analyses in a Bayesian analytical framework using a Bayesian generalized linear mixed models implemented in the brms R package (Bürkner, 2017). The predictor variable will be 'set size' (factor of 3 levels: set size 2, set size 4, set size 6), and the covariates gender (factor of 2 levels: male, female), EHI (factor of 2 levels: right-handed, left-handed), and site (factor of 10 levels) will be added. Subsequently, the credible intervals (CIs) of the posterior distributions were calculated from the newly estimated levels of significance. We opted not to calculate Bayes factors for point estimates to determine whether the effect was zero or unequal to zero. This decision was made because these Bayes factors, which rely on the Savage-Dickey ratio, heavily depend on the arbitrary selection of the prior distribution for each effect. Instead, we employed a different approach: we considered a model parameter to be significant if its 95% confidence interval (CI) does not include zero. As suggested by Gelman et al. (2007), the predictors and outcome variables will be scaled to achieve a mean of 0 and a standard deviation of 0.5. For initial prior distributions, uninformative Cauchy priors will be set to a mean of 0 and a standard deviation of 2.5. Please note that the BayesFactor R package, which would provide an ANOVA design, does not provide the capability to specify the precise location of the prior, thus making it unsuitable for implementing a Bayesian sequential updating approach. This approach is necessary for accumulating evidence across various datasets and determining the success of replication.

**Statistical analysis for Hypothesis #2**

In the original paper, VWM capacity was positively correlated with CDA amplitude increase from set size 2 to 4 (when the smaller set size is below typical adult working memory capacity

estimates), but not from set size 4 to 6 (when both set sizes are at or exceed capacity estimates for typical adults). To replicate this, we will calculate the mean CDA amplitude increase from set size 2 to 4, and from set size 4 to 6, for each subject individually. The VWM capacity will be calculated using the same formula as in the original study. This formula was introduced by Pashler (1988) and refined by Cowan (2001). It is based on the assumption that if a person can retain K items from an S-item array, then the changed item should be among the K items being held in memory on (K/S) trials, leading to correct answers on (K/S) trials where an item changed. The formula accounts for the false alarm rate to adjust for guessing and is expressed as $K = S \times (H - F)$, where K is the memory capacity, S is the set size, H is the observed hit rate in the given set size, and F is the false alarm rate in the given set size. The resulting K scores from all set sizes (i.e., 2, 4, 6) will be used to compute an average K score, which we will use as the behavioral measure of VWM capacity. The relationship between VWM capacity and an increase in CDA amplitude from two to four items will be statistically tested using Pearson's correlation. The significance level will be set to $p < 0.02$ (one-sided). If the Pearson's correlation reveals a significant positive relationship between VWM capacity and the CDA amplitude increase from two to four items, it will support the hypothesis [H2.1]. Furthermore, we will conduct a Bayesian linear mixed model with a prior assuming the reported correlation coefficient from the original study ($r = .78$). The covariates will include gender, EHI, and site. Again, significance is considered if the 95% confidence interval (CI) of the model parameter does not include zero.

**Replication Success**

Replication success will be assessed for each hypothesis separately and is defined operationally as a statistically significant random-effects meta-analytic estimate (at $p < .02$) combining the results from the different laboratories (with a site as a random effect), in the same direction as in the original study.

Hypothesis [H. 1.3] and [H.2.2] will be analyzed using an equivalence test for meta-analyses (Lakens, 2017). The equivalence test assesses whether the difference in CDA amplitude between arrays of four items and six items is as extreme as the smallest effect size of interest

(SESOI) using the two one-sided tests (TOST) procedure implemented in the R package TOSTER (Caldwell, 2022; Lakens, 2017). To perform TOST, the SESOI and its lower and upper equivalence bounds must be established. For replication studies Simonsohn recommended specifying the equivalence bounds for replication studies using the "small telescopes approach" (Simonsohn, 2015). The idea is to consider the effect size that would give the original study 33% power. If the original study had 33% power, the probability of observing a significant effect, if there was a true effect, is too low to reliably distinguish signal from noise (Simonsohn, 2015). Using the small telescopes approach for hypothesis [H 1.3], the SESOI is d = 0.36. An alternative approach would be to calculate the smallest effect size that has enough power to be detected based on the given sample size and alpha level. With this approach the smallest effect size would be very similar to the small telescope approach (i.e., d = 0.44). Therefore, we decided to define the SESOI based on the "small telescopes approach" (i.e., d = 0.36) as this approach was specifically recommended for replication studies. The TOST procedure is then conducted against these bounds based on the SESOI. If the 90% confidence interval of the meta-analytic effect size falls within the equivalence bounds, the observed meta-analytic effect is statistically equivalent (Lakens, 2017). In order to test hypothesis [2.2.], which postulates that there is no correlation between the subject's VWM capacity and the CDA amplitude increase from four to six items, we will conduct another equivalence test. Similar to hypothesis [1.3.], we will use the small telescope approach to specify the SESOI (i.e., r = ±0.29).

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

## Sensitivity analyses

### Measure of VWM capacity

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

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

# Pilot data

We conducted a pilot study to test the feasibility of proposed methods, to implement sanity checks, and to prepare analysis code for the main analyses. First, we demonstrated the feasibility of the estimation of the individual's amplitude threshold for 1° horizontal eye movement, as described in the exclusion criteria section. To determine the individual participant's exclusion amplitude threshold, we have conducted a separate horizontal EOG saccade calibration task. Participants conducted saccades to the target location which appeared either 3° or 6° away from the fixation cross in the center. The pilot study indicates that the estimation of individual's amplitude threshold can be estimated (Figure 2).



*Figure 2. A: The mean amplitudes for the left (red) and right (blue) saccades are displayed for an individual subject. The black lines indicate the single trials. B: The estimation of the amplitude threshold for 1° visual angle for an individual subject is illustrated.*

Subsequently, the pilot study investigated the outcome-neutral tests as well as set size effects in the color contrast change detection task. Because of the small sample size (N = 3), the pilot dataset has insufficient power (see Power Analysis section) to conduct robust statistical analyses. Therefore, the pilot data results are reported here as descriptive results and presented exclusively for illustrative purposes. We refrain from any interpretation based on these data, which should be based on the sample of the actual study. The data acquisition, preprocessing, and analysis parameters were identical to those described for the planned study. Pilot data analysis showed the feasibility of the proposed analysis in the main study. First, we demonstrated the presence of an asymmetry between contra- or ipsilateral electrode

clusters time-locked to the memory array (i.e., outcome-neutral test) (see Figure 3A). In addition, pilot data indicate descriptive differences in CDA amplitude between set size 2 and set size 4 and between set size 2 and 6 respectively (Figure 3B).
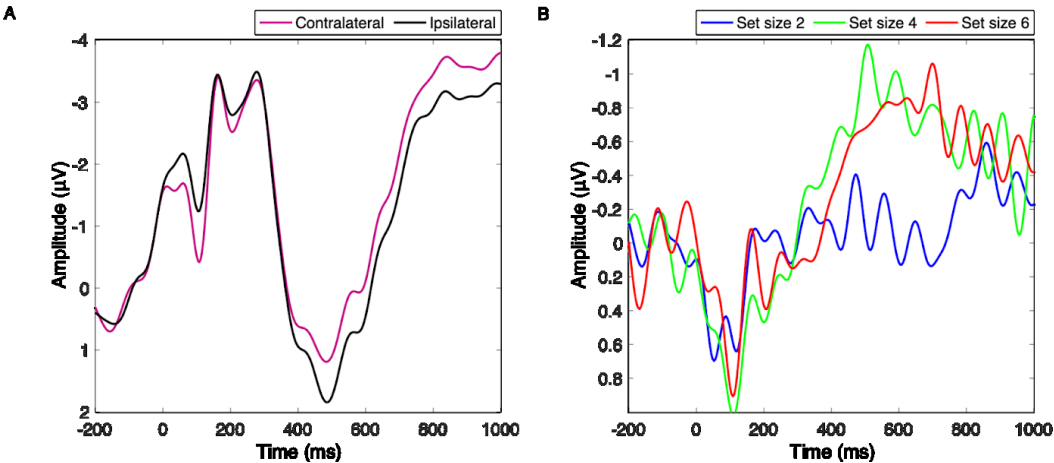


*Figure 3: A: Grand averaged ERP waveforms time-locked to the memory array averaged across the lateral occipital and posterior parietal electrode sites. Note that, by convention, negative voltage is plotted upwards. B: lateralized ERPs are plotted for set size 2, 4 and 6.*

# Acknowledgements

# Funding

# Competing Interests

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

# References

Adam, K. C. S., Robison, M. K., & Vogel, E. K. (2018). Contralateral Delay Activity Tracks Fluctuations in Working Memory Performance. *Journal of Cognitive Neuroscience*, *30*(9), 1229–1240.

Asp, I. E., Störmer, V. S., & Brady, T. F. (2021). Greater Visual Working Memory Capacity for Visually Matched Stimuli When They Are Perceived as Meaningful. *Journal of Cognitive Neuroscience*, 1–17.

Baldwin, S. A., Larson, M. J., & Clayson, P. E. (2015). The dependability of electrophysiological measurements of performance monitoring in a clinical sample: A generalizability and decision analysis of the ERN and Pe. *Psychophysiology*, *52*(6), 790–800.

Bartoš, F., Maier, M., Shanks, D. R., Stanley, T. D., Sladekova, M., & Wagenmakers, E.-J. (2023). Meta-analyses in psychology often overestimate evidence for and size of effects. Royal Society Open Science, 10(7), 230224.

Bigdely-Shamlo, N., Touryan, J., Ojeda, A., Kothe, C., Mullen, T., & Robbins, K. (2020). Automated EEG mega-analysis I: Spectral and amplitude characteristics across studies. *NeuroImage*, *207*, 116361.

Bleuzé, A., Mattout, J., & Congedo, M. (2021). Transfer Learning for the Riemannian Tangent Space: Applications to Brain-Computer Interfaces. *2021 International Conference on Engineering and Emerging Technologies (ICEET)*, 1–6.

Brady, T. F., Robinson, M. M., Williams, J. R., & Wixted, J. T. (2022). Measuring memory is harder than you think: How to avoid problematic measurement practices in memory research. *Psychonomic Bulletin & Review*. https://doi.org/10.3758/s13423-022-02179-w

Brady, T. F., Störmer, V. S., & Alvarez, G. A. (2016). Working memory is not fixed-capacity: More active storage capacity for real-world objects than for simple stimuli. *Proceedings of the National Academy of Sciences of the United States of America*, *113*(27), 7459–7464.

Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, *10*(4), 433–436.

Brennan, R. L. (2005). Generalizability theory. *Educational Measurement Issues and Practice*, *11*(4), 27–34.

Brisson, B., & Jolicoeur, P. (2007). A psychological refractory period in access to visual short-term memory and the deployment of visual-spatial attention: multitasking processing deficits revealed by event-related potentials. *Psychophysiology*, *44*(2), 323–333.

Caldwell, A. R. (2022). *Exploring Equivalence Testing with the Updated TOSTER R Package*. https://doi.org/10.31234/osf.io/ty8de

Champely, S. (2020). *Basic Functions for Power Analysis [R package pwr version 1.3-0]*. https://CRAN.R-project.org/package=pwr

Clayson, P. E., Baldwin, S. A., & Larson, M. J. (2021). Evaluating the internal consistency of subtraction-based and residualized difference scores: Considerations for psychometric reliability analyses of event-related potentials. *Psychophysiology*, *58*(4), e13762.

Clayson, P. E., Brush, C. J., & Hajcak, G. (2021). Data quality and reliability metrics for event-related potentials (ERPs): The utility of subject-level reliability. *International Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *165*, 121–136.

Clayson, P. E., Carbine, K. A., Baldwin, S. A., & Larson, M. J. (2019). Methodological reporting behavior, sample sizes, and statistical power in studies of event-related potentials: Barriers to reproducibility and replicability. *Psychophysiology*, *56*(11), e13437.

Clayson, P. E., Carbine, K. A., Baldwin, S. A., Olsen, J. A., & Larson, M. J. (2021). Using generalizability theory and the ERP Reliability Analysis (ERA) Toolbox for assessing test-retest reliability of ERP scores part 1: Algorithms, framework, and implementation. *International Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *166*, 174–187.

Clayson, P. E., & Miller, G. A. (2017a). ERP Reliability Analysis (ERA) Toolbox: An open-source toolbox for analyzing the reliability of event-related brain potentials. *International Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *111*, 68–79.

Clayson, P. E., & Miller, G. A. (2017b). Psychometric considerations in the measurement of event-related brain potentials: Guidelines for measurement and reporting. *International*

*Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *111*, 57–67.

Cowan, N. (2001). The magical number 4 in short-term memory: a reconsideration of mental storage capacity. *The Behavioral and Brain Sciences*, *24*(1), 87–114; discussion 114–185.

Cowan, N., Elliott, E. M., Scott Saults, J., Morey, C. C., Mattox, S., Hismjatullina, A., & Conway, A. R. A. (2005). On the capacity of attention: its estimation and its role in working memory and cognitive aptitudes. *Cognitive Psychology*, *51*(1), 42–100.

de Cheveigné, A. (2020). ZapLine: A simple and effective method to remove power line artifacts. *NeuroImage*, *207*, 116356.

Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *Journal of Neuroscience Methods*, *134*(1), 9–21.

Emrich, S. M., Al-Aidroos, N., Pratt, J., & Ferber, S. (2009). Visual search elicits the electrophysiological marker of visual working memory. *PloS One*, *4*(11), e8042.

Feldmann-Wüstefeld, T. (2021). Neural measures of working memory in a bilateral change detection task. *Psychophysiology*, *58*(1), e13683.

Feldmann-Wüstefeld, T., Vogel, E. K., & Awh, E. (2018). Contralateral Delay Activity Indexes Working Memory Storage, Not the Current Focus of Spatial Attention. *Journal of Cognitive Neuroscience*, *30*(8), 1185–1196.

Feuerstahler, L. M., Luck, S. J., MacDonald, A., 3rd, & Waller, N. G. (2019). A note on the identification of change detection task models to measure storage capacity and attention in visual working memory. *Behavior Research Methods*, *51*(3), 1360–1370.

Forsberg, A., Adams, E. J., & Cowan, N. (2022). Why does visual working memory ability improve with age: More objects, more feature detail, or both? A registered report. *Developmental Science*, e13283.

Garrett-Ruffin, S., Hindash, A. C., Kaczkurkin, A. N., Mears, R. P., Morales, S., Paul, K., Pavlov, Y. G., & Keil, A. (2021). Open science in psychophysiology: An overview of challenges and

emerging solutions. *International Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *162*, 69–78.

Gelman, A., & Loken, E. (2013). The garden of forking paths: Why multiple comparisons can be a problem, even when there is no "fishing expedition" or "p-hacking" and the research hypothesis was posited ahead of time. *Department of Statistics, Columbia University, 348*, 1–17.

Hakim, N., Adam, K. C. S., Gunseli, E., Awh, E., & Vogel, E. K. (2019). Dissecting the Neural Focus of Attention Reveals Distinct Processes for Spatial Attention and Object-Based Storage in Visual Working Memory. *Psychological Science*, *30*(4), 526–540.

Heuer, A., & Schubö, A. (2016). The Focus of Attention in Visual Working Memory: Protection of Focused Representations and Its Individual Variation. *PloS One, 11*(4), e0154228.

Jennings, J. R., & Wood, C. C. (1976). Letter: The epsilon-adjustment procedure for repeated-measures analyses of variance. *Psychophysiology*, *13*(3), 277–278.

Jongbloed-Pereboom, M., Nijhuis-van der Sanden, M. W. G., & Steenbergen, B. (2019). Explicit and implicit motor sequence learning in children and adults; the role of age and visual working memory. *Human Movement Science*, *64*, 1–11.

Klaver, P., Talsma, D., Wijers, A. A., Heinze, H. J., & Mulder, G. (1999). An event-related brain potential correlate of visual short-term memory. *Neuroreport*, *10*(10), 2001–2005.

Lakens, D. (2017). Equivalence Tests: A Practical Primer for t Tests, Correlations, and Meta-Analyses. *Social Psychological and Personality Science*, *8*(4), 355–362.

Lane, D. M., & Dunlap, W. P. (1978). Estimating effect size: Bias resulting from the significance criterion in editorial decisions. The British Journal of Mathematical and Statistical Psychology, 31(2), 107–112.

Levine, T. R., Asada, K. J., & Carpenter, C. (2009). Sample Sizes and Effect Sizes are Negatively Correlated in Meta-Analyses: Evidence and Implications of a Publication Bias Against NonSignificant Findings. Communication Monographs, 76(3), 286–302.

Liesefeld, H. R., & Müller, H. J. (2019). Current directions in visual working memory research: An introduction and emerging insights. *British Journal of Psychology* , *110*(2), 193–206.

Lin, L., & Chu, H. (2018). Quantifying publication bias in meta-analysis. Biometrics, 74(3), 785–794.

Lotfi, S., Ward, R., Mathew, A., Shokoohi-Yekta, M., Rostami, R., Motamed-Yeganeh, N., Christine, C., & Lee, H.-J. (2022). Limited visual working memory capacity in children with dyslexia: An ERP study. *NeuroRegulation*, *9*(2), 98–109.

Luck, S. J. (2014). *An Introduction to the Event-Related Potential Technique, second edition*. MIT Press.

Luck, S. J., & Gaspelin, N. (2017). How to get statistically significant effects in any ERP experiment (and why you shouldn't). *Psychophysiology*, *54*(1), 146–157.

Luck, S. J., Stewart, A. X., Simmons, A. M., & Rhemtulla, M. (2021). Standardized measurement error: A universal metric of data quality for averaged event-related potentials. *Psychophysiology*, *58*(6), e13793.

Luck, S. J., & Vogel, E. K. (2013). Visual working memory capacity: from psychophysics and neurobiology to individual differences. *Trends in Cognitive Sciences*, *17*(8), 391–400.

Luria, R., Balaban, H., Awh, E., & Vogel, E. K. (2016). The contralateral delay activity as a neural measure of visual working memory. *Neuroscience and Biobehavioral Reviews*, *62*, 100–108.

Maman, G., Yair, O., Eytan, D., & Talmon, R. (2019). Domain Adaptation Using Riemannian Geometry of Spd Matrices. *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 4464–4468.

Mellot, A., Collas, A., Rodrigues, P. L. C., Engemann, D., & Gramfort, A. (2023). Harmonizing and aligning M/EEG datasets with covariance-based techniques to enhance predictive regression modeling. In *bioRxiv* (p. 2023.04.27.538550). https://doi.org/10.1101/2023.04.27.538550

Miyake, A., & Shah, P. (Eds.). (1999). *Models of working memory: Mechanisms of active maintenance and executive control*. *506*. https://doi.org/10.1017/CBO9781139174909

Naveh-Benjamin, M., & Cowan, N. (2023). The roles of attention, executive function and knowledge in cognitive ageing of working memory. *Nature Reviews Psychology*, 1–15.

Ngiam, W. X. Q., Adam, K. C. S., Quirk, C., Vogel, E. K., & Awh, E. (2021). Estimating the statistical power to detect set-size effects in contralateral delay activity. *Psychophysiology*, *58*(5), e13791.

Oldfield, R. C. (1971). The assessment and analysis of handedness: the Edinburgh inventory. *Neuropsychologia*, *9*(1), 97–113.

Pashler, H. (1988). Familiarity and visual change detection. *Perception & Psychophysics*, *44*(4), 369–378.

Paul, M., Govaart, G. H., & Schettino, A. (2021). Making ERP research more transparent: Guidelines for preregistration. *International Journal of Psychophysiology: Official Journal of the International Organization of Psychophysiology*, *164*, 52–63.

Pavlov, Y. G., Adamian, N., Appelhoff, S., Arvaneh, M., Benwell, C. S. Y., Beste, C., Bland, A. R., Bradford, D. E., Bublatzky, F., Busch, N. A., Clayson, P. E., Cruse, D., Czeszumski, A., Dreber, A., Dumas, G., Ehinger, B., Ganis, G., He, X., Hinojosa, J. A., … Mushtaq, F. (2021). #EEGManyLabs: Investigating the replicability of influential EEG experiments. *Cortex; a Journal Devoted to the Study of the Nervous System and Behavior*, *144*, 213–229.

Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: transforming numbers into movies. *Spatial Vision*, *10*(4), 437–442.

Perron, R., Lefebvre, C., Robitaille, N., Brisson, B., Gosselin, F., Arguin, M., & Jolicoeur, P. (2009). Attentional and anatomical considerations for the representation of simple stimuli in visual short-term memory: evidence from human electrophysiology. *Psychological Research*, *73*(2), 222–232.

Pickering, S. J. (2006). Working memory in dyslexia. In T. P. Alloway (Ed.), *Working memory and neurodevelopmental disorders , (pp* (Vol. 306, pp. 7–40). Psychology Press, xiii.

Quirk, C., Adam, K. C. S., & Vogel, E. K. (2020). No Evidence for an Object Working Memory Capacity Benefit with Extended Viewing Time. *eNeuro*, *7*(5). https://doi.org/10.1523/ENEURO.0150-20.2020

Rodrigues, P. L. C., Jutten, C., & Congedo, M. (2019). Riemannian Procrustes Analysis: Transfer Learning for Brain–Computer Interfaces. *IEEE Transactions on Biomedical Engineering*, *66*(8), 2390–2401.

Rothstein, H., Sutton, A. J., & Borenstein, M. (Eds.). (2005). Publication bias in meta-analysis: Prevention, assessment and adjustments [PDF]. Wiley-Blackwell.

Roy, Y., & Faubert, J. (2023). Is the Contralateral Delay Activity (CDA) a robust neural correlate for Visual Working Memory (VWM) tasks? A reproducibility study. *Psychophysiology*, *60*(2), e14180.

Schneider, D., Barth, A., Getzmann, S., & Wascher, E. (2017). On the neural mechanisms underlying the protective function of retroactive cuing against perceptual interference: Evidence by event-related potentials of the EEG. *Biological Psychology*, *124*, 47–56.

Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations stabilize? In *Journal of Research in Personality* (Vol. 47, Issue 5, pp. 609–612). https://doi.org/10.1016/j.jrp.2013.05.009

Simonsohn, U. (2015). Small telescopes: detectability and the evaluation of replication results. *Psychological Science*, *26*(5), 559–569.

Sundre, D. L. (1993). Book Reviews : Generalizability Theory: A Primer, by Richard J. Shavelson and Noreen M. Webb. Newbury Park, CA: Sage Publications, 1991,137 pp. In *Evaluation Practice* (Vol. 14, Issue 2, pp. 207–209). https://doi.org/10.1177/109821409301400219

Unsworth, N., Fukuda, K., Awh, E., & Vogel, E. K. (2015). Working memory delay activity predicts individual differences in cognitive abilities. *Journal of Cognitive Neuroscience*, *27*(5), 853–865.

Villena-González, M., Rubio-Venegas, I., & López, V. (2020). Data from brain activity during visual working memory replicates the correlation between contralateral delay activity and memory capacity. *Data in Brief*, *28*, 105042.

Vogel, E. K., & Machizawa, M. G. (2004). Neural activity predicts individual differences in visual working memory capacity. *Nature*, *428*(6984), 748–751.

von Bastian, C. C., Belleville, S., Udale, R. C., Reinhartz, A., Essounni, M., & Strobach, T. (2022). Mechanisms underlying training-induced cognitive change. *Nature Reviews Psychology*, *1*(1), 30–41.

Wang, J., Huo, S., Wu, K. C., Mo, J., Wong, W. L., & Maurer, U. (2022). Behavioral and neurophysiological aspects of working memory impairment in children with dyslexia. *Scientific Reports*, *12*(1), 12571.

Westheimer, G. (1954a). Mechanism of saccadic eye movements. *A.M.A. Archives of Ophthalmology*, *52*(5), 710–724.

Westheimer, G. (1954b). Eye movement responses to a horizontally moving visual stimulus. *A.M.A. Archives of Ophthalmology*, *52*(6), 932–941.

Widmann, A., & Schröger, E. (2012). Filter effects and filter artifacts in the analysis of electrophysiological data. *Frontiers in Psychology*, *3*, 233.

Williams, J. R., Robinson, M. M., Schurgin, M. W., Wixted, J. T., & Brady, T. F. (2022). You cannot "count" how many items people remember in visual working memory: The importance of signal detection–based measures for understanding change detection performance. In *Journal of Experimental Psychology: Human Perception and Performance* (Vol. 48, Issue 12, pp. 1390–1409). https://doi.org/10.1037/xhp0001055

Yang, Y., Sánchez-Tójar, A., O'Dea, R. E., Noble, D. W. A., Koricheva, J., Jennions, M. D., Parker, T. H., Lagisz, M., & Nakagawa, S. (2023). Publication bias impacts on effect size, statistical power, and magnitude (Type M) and sign (Type S) errors in ecology and evolutionary biology. BMC Biology, 21(1), 71.