

University of Groningen

## Dynamic hidden states underlying working-memory-guided behavior

Wolff, Michael J.; Jochim, Janina; Akyürek, Elkan G.; Stokes, Mark G.

*Published in:*  
Nature neuroscience

*DOI:*  
[10.1038/nn.4546](https://doi.org/10.1038/nn.4546)

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*  
Publisher's PDF, also known as Version of record

*Publication date:*  
2017

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Wolff, M. J., Jochim, J., Akyürek, E. G., & Stokes, M. G. (2017). Dynamic hidden states underlying working-memory-guided behavior. *Nature neuroscience*, 20(6), 864-871. <https://doi.org/10.1038/nn.4546>

### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

*Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.*

# Dynamic hidden states underlying working-memory-guided behavior

Michael J Wolff<sup>1,2</sup>, Janina Jochim<sup>1</sup>, Elkan G Akyürek<sup>2</sup> & Mark G Stokes<sup>1</sup>

Recent theoretical models propose that working memory is mediated by rapid transitions in ‘activity-silent’ neural states (for example, short-term synaptic plasticity). According to the dynamic coding framework, such hidden state transitions flexibly configure memory networks for memory-guided behavior and dissolve them equally fast to allow forgetting. We developed a perturbation approach to measure mnemonic hidden states in an electroencephalogram. By ‘pinging’ the brain during maintenance, we show that memory-item-specific information is decodable from the impulse response, even in the absence of attention and lingering delay activity. Moreover, hidden memories are remarkably flexible: an instruction cue that directs people to forget one item is sufficient to wipe the corresponding trace from the hidden state. In contrast, temporarily unattended items remain robustly coded in the hidden state, decoupling attentional focus from cue-directed forgetting. Finally, the strength of hidden-state coding predicts the accuracy of working-memory-guided behavior, including memory precision.

Working memory (WM) is a core cognitive function critical for flexible, intelligent behavior<sup>1</sup>. Until recently, it was widely assumed that information is maintained in WM by maintaining specific activity states that represent the specific memoranda<sup>2,3</sup>. However, accumulating evidence increasingly shows that successful maintenance in WM is not strictly dependent on an unbroken chain of corresponding delay activity<sup>4</sup> and that item-specific activity states could reflect other cognitive processes. For example, in monkey studies, persistent activity ramps up with expectation of the probe<sup>5–8</sup>. Similarly, in humans, it has been shown that unattended WM content is not reflected in the neural signal, even when it is still clearly maintained<sup>9–11</sup>. Evidence for WM in the absence of persistent delay activity suggests that WM can be maintained in activity-silent neural states<sup>4</sup>.

Recent theories acknowledge that brain activity is highly dynamic, even when the contents of working memory remain stable<sup>12</sup>. Multiple neurophysiological mechanisms could underlie such dynamics<sup>13–15</sup>. According to a dynamic coding model of WM<sup>4</sup>, behaviorally relevant sensory input drives a memory-item-specific neural response, which triggers an item-specific change in the functional state of the system. Depending on the precise neural mechanism, this functional state could be activity-silent (for example, short-term synaptic plasticity<sup>14,16–19</sup>) and maintained throughout the memory delay to serve as the neural context for subsequent processing. Items in WM would be read out via the context-dependent response to a probe stimulus during recall<sup>13,20</sup>. Crucially, this model predicts that dynamic hidden states are constructed when new information is encoded and dissolved as soon as it is forgotten. This model also predicts that dynamic hidden states should determine the quality of a representation maintained in WM.

To probe hidden neural states, we developed a functional perturbation approach to ‘ping’ the brain. Analogously to the use of active

sonar (or echolocation), the response to a well-characterized impulse stimulus can be used to infer the current state of the system<sup>4,13</sup>. We recently validated this general approach using noninvasive electroencephalography (EEG) in a proof-of-principle study<sup>21</sup>. The presentation of a high-contrast neutral visual stimulus evoked neural activity that clearly discriminated the previously presented visual stimulus. Here we exploit this approach to track the functional dynamics of hidden states for WM.

Across two experiments, we showed that the content of WM could be decoded from the impulse response during the maintenance interval, while forgotten information left effectively no trace. In Experiment 2, we demonstrated robust hidden-state representations for unattended content in WM, providing a plausible mechanism for maintenance that is independent of the activity associated with the focus of attention. Finally, we also found evidence that the quality of working memory varied with the decodability of these hidden states.

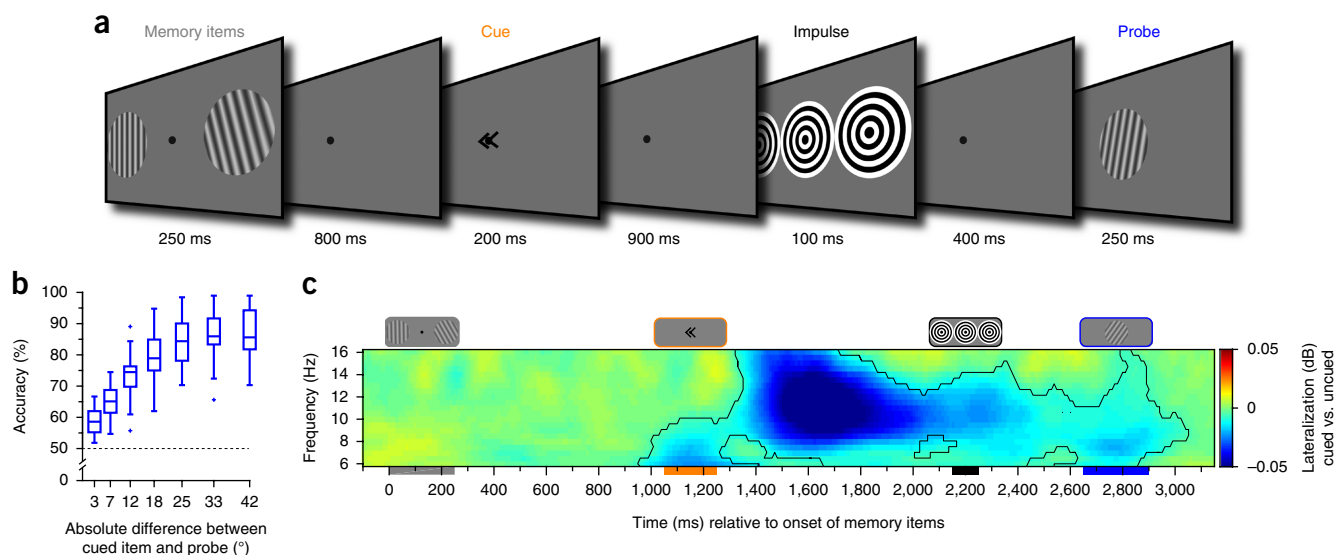
## RESULTS

### Experiment 1

In Experiment 1, 30 human participants performed a visual WM task while EEG was recorded. At the beginning of each trial (**Fig. 1a**), two memory items were presented, but a retrospective cue (retro-cue) presented during the delay instructed participants which item would actually be probed<sup>22,23</sup>. The other item could be simply forgotten. The retro-cue in this design was essential to differentiate WM from basic stimulation history<sup>24</sup>. During a subsequent memory delay, we then presented a high-contrast ‘impulse’ stimulus. Memory performance for the cued item was tested after the impulse by a centrally presented memory probe (**Fig. 1b**). Time–frequency decomposition of lateralized activity in posterior sensors (**Fig. 1c**) showed significant lateralization

<sup>1</sup>Department of Experimental Psychology, University of Oxford, Oxford, UK. <sup>2</sup>Department of Experimental Psychology, University of Groningen, Groningen, the Netherlands. Correspondence should be addressed to M.G.S. ([mark.stokes@psy.ox.ac.uk](mailto:mark.stokes@psy.ox.ac.uk)).

Received 26 January; accepted 15 March; published online 17 April 2017; doi:10.1038/nn.4546



**Figure 1** Experiment 1 task structure, behavioral performance and attention-related alpha-band activity. **(a)** Trial schematic. Two memory items were presented (randomly oriented grating stimuli), and participants were instructed to memorize both orientations. A retro-cue then indicated which item would actually be tested at the end of the current trial (100% valid). The impulse stimulus (high-contrast, task-irrelevant visual input) was then presented during the subsequent delay while participants should have only the cued item in WM. At the end of the trial, a forced-choice probe was presented at the center of the screen. Participants indicated whether the probe was rotated clockwise or anticlockwise relative to the orientation of the cued item. **(b)** Boxplots show WM accuracy as a function of the absolute angular difference (in degrees) between the memory item and the probe. Center line indicates the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5 $\times$  the interquartile range. Extreme values are shown separately (crosses). Dashed line indicates 50% accuracy, chance. **(c)** Time–frequency representation of the difference between the contra- and ipsilateral posterior electrodes relative to the cued hemifield. The highlighted cluster in the alpha-frequency band (8–12 Hz) indicates significant contralateral desynchronization (permutation test,  $n = 30$ , cluster-forming threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Colored bars under the  $x$  axis represent the timings of the corresponding stimuli illustrated on top.

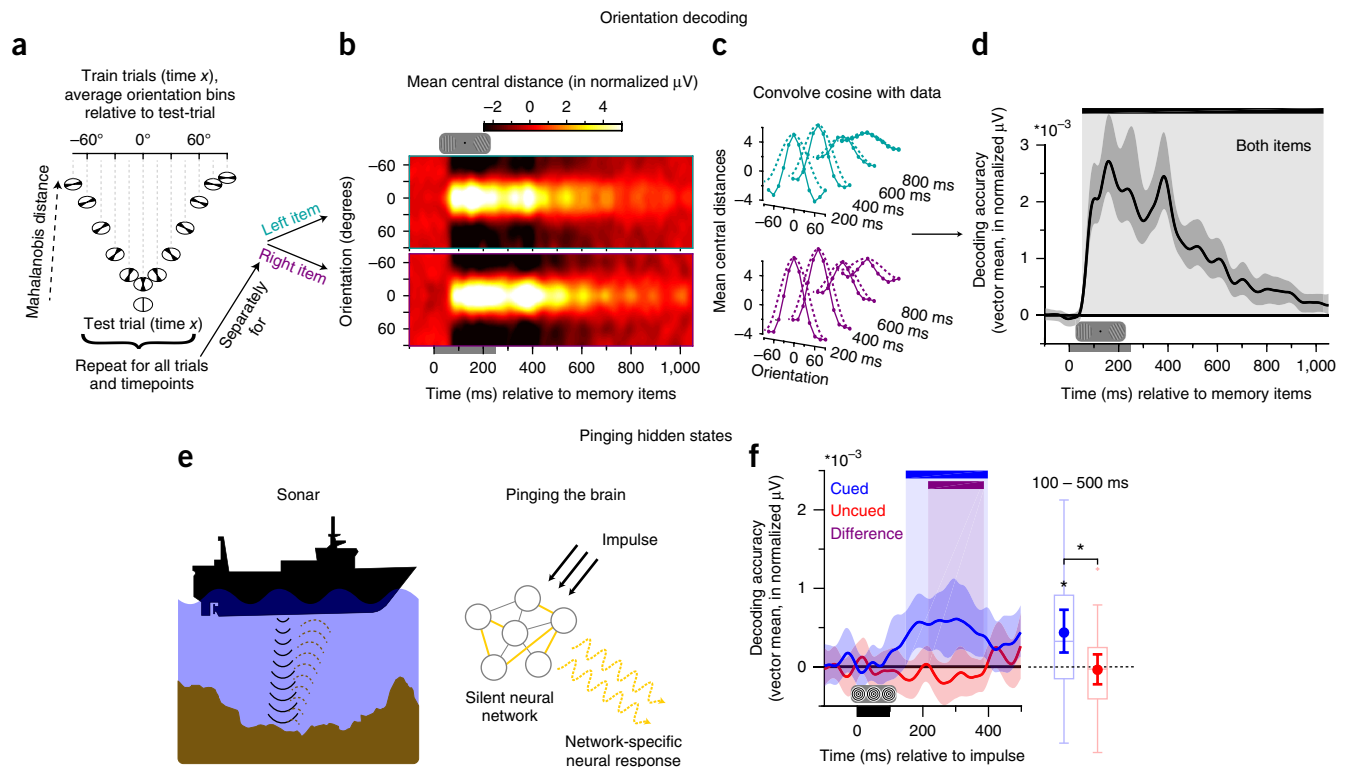
in the alpha range (8–12 Hz) after the presentation of the cue (permutation test,  $n = 30$ ,  $P < 0.001$ , corrected; cluster-forming threshold  $P < 0.05$ ). This pattern was consistent with a shift in spatial attention<sup>25</sup> according to the retro-cue, which confirmed that the cue manipulation was effective.

**Decoding parametric memory items.** To decode the memory items used in this experiment, we developed a parametric variant of distance-based discrimination (Fig. 2a–d and Online Methods). As shown in Figure 2a, this capitalized on the parametric structure of the stimulus space<sup>26</sup>, whilst maintaining the statistical advantages of the Mahalanobis distance metric used in previous EEG and magnetoencephalography decoding studies<sup>21,27</sup> (Online Methods). To summarize briefly here: for a given trial, we compared the activity pattern across electrodes to the corresponding activity pattern observed in the remaining trials, averaged by orientation-difference to the test trial (at a bin width of 30°). This procedure was repeated for all trials and all timepoints. If the pattern of activity contained information about item orientation, we expect greater pattern dissimilarity (i.e., Mahalanobis distance) at larger angular differences. Figure 2b shows distance as a function of reference angle and time after the presentation of the left and right items separately. Distance values were then converted into a decoding accuracy score (Fig. 2c) and averaged across both items at each timepoint (Fig. 2d). Item orientation could be decoded from 56 ms until 1,026 ms after onset (permutation test,  $n = 30$ ,  $P < 0.001$ , corrected; cluster-forming threshold  $P < 0.05$ ). This was consistent with previous empirical evidence that EEG is sufficiently sensitive to detect subtle differences in scalp-level activity patterns associated with different stimulus orientation<sup>21</sup>. The current decoding results further validated the utility of multivariate pattern analysis for two simultaneously presented orientation gratings.

For completeness, we also decoded item-specific orientation during the retro-cue epoch (Supplementary Fig. 1).

**Pinging hidden states.** On the basis of the dynamic coding framework, we hypothesized that the input–output mapping of neural circuits maintaining information in WM should systematically reflect the memory content<sup>4</sup>. We tested this using an impulse stimulus to ping potentially hidden neural states (Fig. 2e). As predicted, the impulse-specific response clearly differentiated the content of WM (Fig. 2f), even though the driving input (ping) was held constant on each trial. The decodability of the cued item showed a significant cluster from 148 to 398 ms after impulse stimulus onset (permutation test,  $n = 30$ ,  $P = 0.002$ , corrected; cluster-forming threshold  $P < 0.05$ ). Average decodability from 100 to 500 ms was also significant ( $P = 0.004$ ), and cued-item decoding was also higher than task-irrelevant (uncued) item decoding (cluster: 216 to 386 ms,  $P = 0.009$ , corrected; average:  $P = 0.028$ ). Indeed, the uncued item showed no evidence for decoding (no corrected clusters; average:  $P = 0.687$ ), suggesting that content can be rapidly purged from WM when instructed, leaving effectively no trace in the neural state.

To test whether the impulse response reflects a literal ‘reactivation’ of item-specific activity observed during encoding (for example, Fig. 2b), we also examined whether a classifier trained on the activity elicited by the memory stimuli during encoding could be used to decode the memory item during the impulse epoch (and vice versa). However, we found no evidence for significant cross-generalization between discriminative activity patterns during encoding and discriminative activity driven by the impulse (corrected clusters,  $P > 0.347$ ). We propose that the impulse stimulus simply acts as a functional ping to recover hidden states, rather than a literal reactivation of a latent representation<sup>21</sup>.



**Figure 2** Orientation decoding in EEG and pinging hidden states of WM. **(a–d)** Decoding procedure. **(a)** The dissimilarity in the neural pattern between a single trial and all other trials is computed as a function of orientation difference (binned:  $30^\circ$ ). **(b)** Average distance to template of all trials for each timepoint during and after memory item presentation, plotted separately for the left and right memory items (upper and lower, respectively). Distances are mean-centered and sign-reversed (high values are therefore equivalent to a high similarity and short Mahalanobis distance between patterns) for visualization. **(c)** A cosine is convolved with the data. Top, cyan lines show data for the left item; bottom, purple lines show data for the right item. Solid lines are the data. Dashed lines are illustrative for the cosine. **(d)** The vector mean of the convolved tuning curves (i.e., decoding accuracy) over time, averaged over left and right items. Black bar, significant decoding (permutation test,  $n = 30$ , cluster-forming threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Error shading, 95% confidence interval (CI) of the mean. Gray bar on the x axis shows onset of the memory array. **(e)** Pinging hidden states. Analogy to active sonar: differences in hidden state are inferred from differences in the measured response to a well-characterized impulse. **(f)** Decoding results in the impulse epoch. Blue bar, significant decoding of the cued item; purple bar, significant difference in decodability between the cued and uncued item (permutation test,  $n = 30$ , cluster-forming threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Error shading, 95% CI of the mean. Black bar on the x axis shows onset of the impulse stimulus. Right: boxplots and superimposed circles with error bars (mean and 95% CI of the mean) show average decoding from 100 to 500 ms after impulse onset. Data points outside of  $1.5\times$  the interquartile range are shown separately (small crosses). Significant average decoding and significant differences in average decodability between the cued and uncued item are marked by asterisks (permutation test,  $n = 30$ , cued:  $P = 0.004$ ; difference:  $P = 0.028$ ).

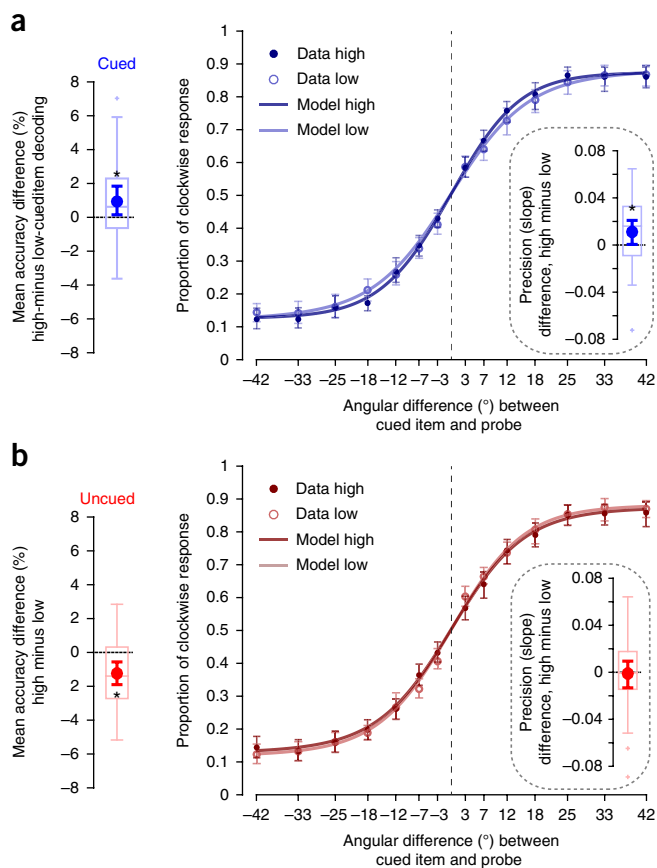
Trial-wise variability in decoding the impulse response also predicted variability in WM performance. Higher-decoding trials of the cued item were accompanied by higher performance than low-decoding trials (permutation test,  $n = 30$ ,  $P = 0.043$ ; **Fig. 3a**). There was also a complementary cost for decoding the uncued item (i.e., a high decoding score for the uncued item led to a decrease in accuracy on the cued item;  $P = 0.002$ ; **Fig. 3b**), suggesting that participants might have failed to discard the uncued item (or simply did not use the cue properly) on some trials, contributing to error in performance. Finally, the difference between the accuracy effect of the cued and the uncued item was also significant (permutation test,  $n = 30$ ,  $P < 0.001$ ).

In principle, the relationship between trial-wise decoding and WM performance may rest on an increase in the guess rate (i.e., due to forgetting or failure to encode), a reduction in precision or both<sup>28,29</sup>. To separate these possible contributions, we modeled the behavioral profile over degrees of angular rotation between the memory item and the probe stimulus (Online Methods; <http://www.palamedestoolbox.org>)<sup>30</sup>. We found that the link to behavior was most likely driven by a decrease in precision (the slope parameter of the model) for weakly

encoded hidden states of WM (permutation test,  $n = 30$ ,  $P = 0.023$ , one-tailed; **Fig. 3a**), while no evidence for an effect in guess rate (the asymptote parameter) was found ( $P = 0.867$ , one-tailed). Modeling the observed uncued item accuracy effect was inconclusive (**Fig. 3b**), with no evidence for either a precision or guess rate effect ( $P = 0.443$  and  $P = 0.184$ , respectively, one-tailed). Finally, we found no evidence that trial-wise item decoding during the initial presentation of the memory stimuli related to memory performance (**Supplementary Fig. 2a**), further suggesting that the relationship between accuracy and decoding triggered by the impulse was not due to a failure to encode the memory item.

## Experiment 2

Recently, it has been proposed that information in WM can be represented in qualitatively different states<sup>31–33</sup>, with attended items encoded in activity states measurable with standard recordings of delay activity, whereas activity-silent states could underlie the representation of currently unattended information in WM. In Experiment 2 ( $n = 19$  subjects) we tested whether unattended but nevertheless remembered information in WM can still be decoded from the



**Figure 3** Relationship between item-specific impulse decoding and WM accuracy. **(a)** Left: difference in overall WM task performance between high- and low-cued item decoding trials. Right: proportion of clockwise responses for high- and low-decoding trials as a function of the angular difference between the memory item and the probe. Error bars, 95% CI of the mean. Inset shows the difference in the slope parameter (a measure of memory precision) between high- and low-decoding trials. Data points outside of 1.5× the interquartile range are shown separately in the boxplots (small crosses). Superimposed circles and error bars, mean and 95% CI of the mean. **(b)** As in **a** but for decoding the uncued item. Significant differences in accuracy or precision between high- and low-decoding trials are highlighted by asterisks (permutation test,  $n = 30$ , cued:  $P = 0.043$  and  $P = 0.023$ , for accuracy and precision, respectively; uncued:  $P < 0.001$  for accuracy). In boxplots, horizontal line indicates the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5× the interquartile range. Extreme values are shown separately (crosses).

impulse response. Again, two memory items were presented at the start of the trial, but both were ultimately relevant as they would both be probed. Priority was manipulated by blocking the order in which items would be probed (Fig. 4a) and instructing participants accordingly. Because there was no other clue as to which item was being probed first or second, nonrandom responses indicated that participants used this blocked information (Fig. 4b). This was further supported by lateralized changes in alpha power (Fig. 4c). During and shortly after the initial presentation of the memory stimuli, there was a relative decrease in power at sensors contralateral to the initially prioritized item, consistent with selective allocation of attention (permutation test,  $n = 19$ ,  $P = 0.023$ , corrected; cluster-forming threshold  $P < 0.05$ ). Moreover, this pattern reversed after the response to the first item ( $P = 0.009$ , corrected), consistent with the assumption that participants then shifted the originally deprioritized item into the focus of attention in WM in preparation for the second probe<sup>34</sup>.

**Decoding during stimulus presentation.** We first analyzed decoding during the initial processing of the memory stimuli. The results were plotted separately as a function of test time (early or late in the trial), as this could be meaningfully classified from the beginning of the trial (Fig. 5a). As expected, decoding the prioritized item (cluster: 74 to 1,200 ms,  $P < 0.001$ , corrected; cluster-forming threshold  $P < 0.05$ ; average:  $P < 0.001$ ), relative to the deprioritized item (cluster: 82 to 542 ms, corrected,  $P < 0.001$ , corrected; average,  $P < 0.001$ ) was more robust (average:  $P = 0.013$ ). While decoding of the unattended item dropped to chance relatively quickly after item presentation, the attended item showed significant decoding until the end of the epoch, replicating previous evidence showing that maintenance of only attended WM items is represented in the recorded brain activity patterns<sup>9–11</sup>.

The difference between attended and unattended item-maintenance in WM was even more apparent when comparing their cross-temporal decoding matrices. Minimal cross-temporal generalization during and shortly after memory item presentation suggested highly dynamic item encoding: orientation-discriminative patterns change over time. This was supported by significant dynamic coding clusters during item encoding for both the early- and late-tested items, where off-diagonal timepoints showed significantly lower decodability than both corresponding on-diagonal timepoints (permutation test,  $n = 19$ , cluster-defining threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ; Fig. 5b and Online Methods). However, the attended item clearly showed a more time-invariant decoding pattern at the end of the epoch than the unattended item, apparent due to both significantly higher decodability on the same timepoint as well as cross-timepoint decoding ( $n = 19$ ,  $P = 0.023$ , corrected; cluster-forming threshold  $P < 0.05$ ; Fig. 5b). This further suggests that while the attended item also had a corresponding WM maintenance signature in stable activity patterns, the unattended item did not.

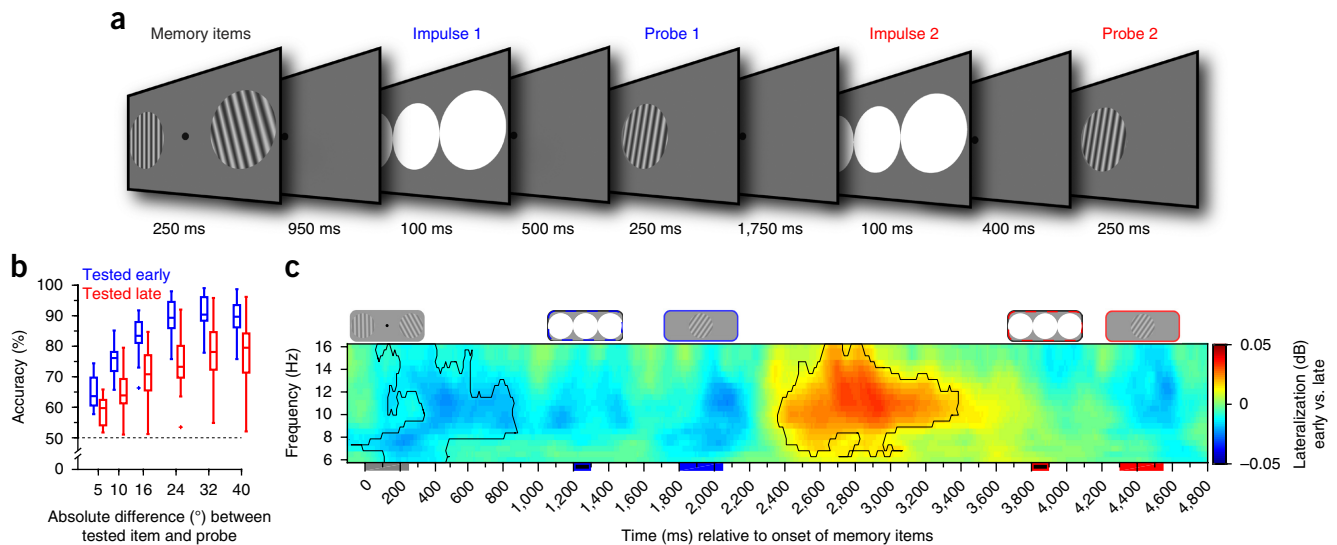
**Decoding of the impulse responses.** Critically, we found that both the attended (clusters: 80 to 308 ms,  $P = 0.004$ ; and 332 ms to 434 ms,  $P = 0.031$ , corrected; average:  $P < 0.001$ ) and unattended items (cluster: 172 to 306 ms,  $P = 0.011$ , corrected; average:  $P = 0.045$ ) were decodable in the first impulse response (Fig. 6a). This contrasted with the clear cueing differences observed in Experiment 1 and suggested that multiple items can be encoded in hidden states and revealed by the impulse, even if only one item is in the focus of attention. It is worth noting, however, that the decodability of the attended item was significantly higher than that of the unattended item (average:  $P = 0.031$ ), consistent with the behavioral evidence for relatively better memory for the initially prioritized item.

We found no evidence for a relationship between trial-wise differences in alpha lateralization and WM item decodability of the impulse response for either the attended or unattended item (Supplementary Fig. 3). This further suggests that the item-specific impulse response does not even vary with trial-wise differences in the focus of attention.

We also found that the remaining relevant and initially unattended item could also be decoded in the second impulse response (cluster: 196 to 326 ms,  $P = 0.016$ , corrected; average:  $P = 0.012$ ), while decoding the initially prioritized item failed to reach significance in this epoch (clusters:  $P > 0.109$ , corrected; average:  $P = 0.112$ ; Fig. 6b). The now-deprioritized item was presumably cleared from the hidden state because it was no longer relevant, similarly to the forgetting observed after the retro-cue from Experiment 1.

Again, we also tested for cross-generalization between the decodable patterns of the memory-items epoch (Fig. 5a) and the impulse epochs (Fig. 6a,b). However, as in Experiment 1, we found no evidence



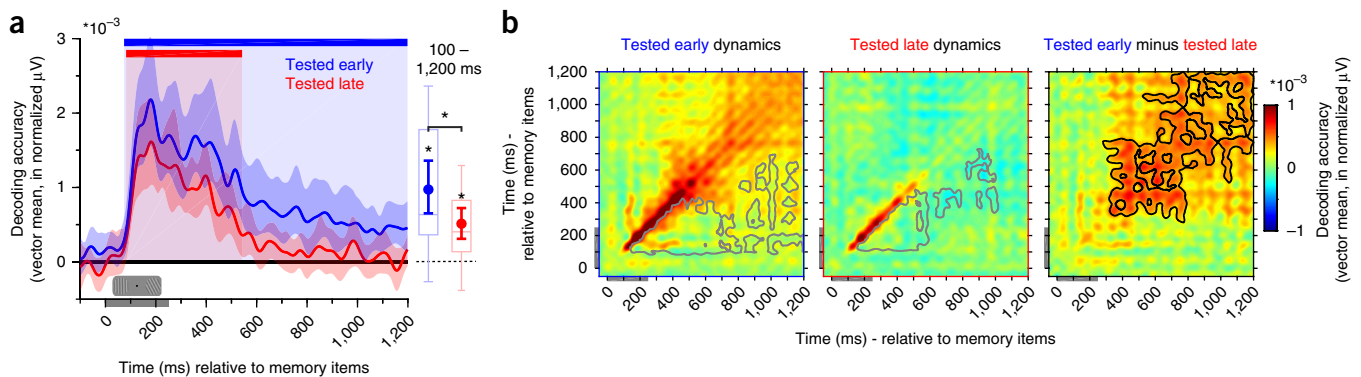


**Figure 4** Experiment 2 task structure, behavioral performance and attention-related alpha-band activity. **(a)** Trial schematic. Two memory items were presented. Participants were instructed to maintain both items and were told at the start of each block the order in which the items would be tested. The first impulse was presented within the first memory delay (maintain both items but attend the prioritized item), after which the prioritized item was probed. The second impulse was presented during the subsequent memory delay (maintain and attend only the now-prioritized item), after which the remaining item was probed. **(b)** Boxplots show the accuracy of the early- and late-tested items as a function of the absolute angular difference (in degrees) between the memory item and the probe. In boxplots, horizontal line indicates the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5x the interquartile range. Extreme values are shown separately (crosses). **(c)** Time–frequency representation of the difference between the contra- and ipsilateral posterior electrodes relative to the presentation side of the early-tested memory items. Highlighted areas indicate significant difference (permutation test,  $n = 19$ , cluster-forming threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ).

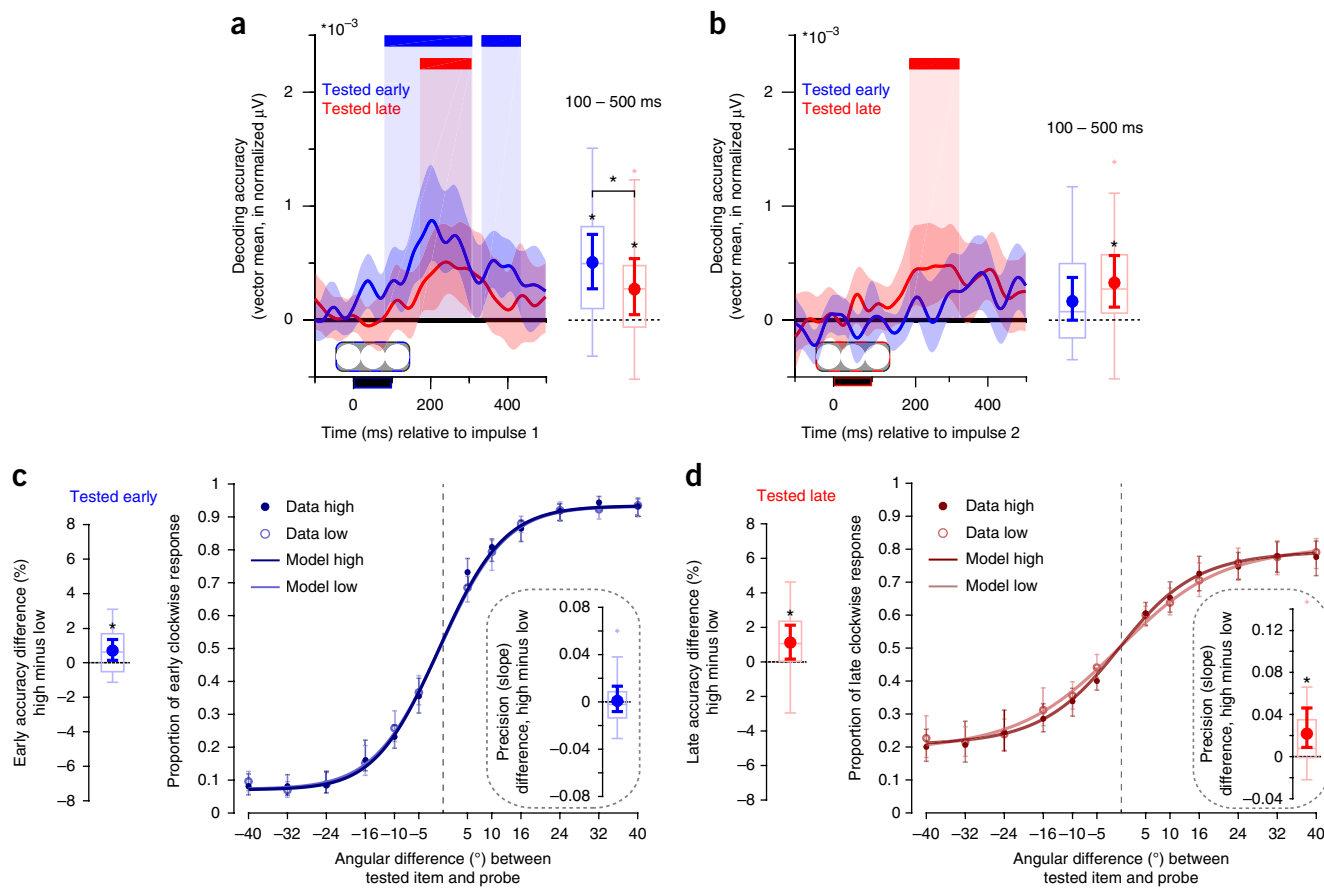
that the impulse literally reactivated activity patterns associated with initial encoding for either item (all corrected clusters:  $P > 0.32$ ).

There was also a positive relationship between trial-wise decoding of the attended items at the first and at the second impulse with WM performance (early:  $P = 0.038$ , Fig. 6c; late:  $P = 0.04$ ; Fig. 6d), replicating and extending the findings of Experiment 1. As in Experiment 1, we modeled the behavioral profile to test whether the positive relationship between decoding and task performance was due to an

increase in precision and/or a decrease in the guess rate. While the modeling results were inconclusive for the early-tested item (precision:  $P = 0.399$ , one-tailed; guess rate:  $P = 0.329$ , one-tailed; Fig. 6c), there was evidence for an effect of WM precision for the late item (precision:  $P = 0.006$ , one-tailed; guess rate:  $P = 0.942$ , one-tailed; Fig. 6d), replicating the precision effect of Experiment 1. Note that there was again no relationship between accuracy and item decoding during the encoding phase (Supplementary Fig. 2b).



**Figure 5** Priority-dependent encoding and maintenance in WM. **(a)** Decodability of the early-test item (blue) and the late-tested item (red) during memory item presentation. Blue and red bars, significant decoding clusters for the early- and late-tested items, respectively (permutation test,  $n = 19$ , cluster-defining threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Error shading, 95% CI of the mean. Boxplots and superimposed circles with error bars (mean and 95% CI of the mean) represent average decodability from 100 ms after stimulus onset until the end of the epoch. In boxplots, horizontal line indicates the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5x the interquartile range. Significant average decoding and average differences between the decodability of the early and late items are marked by asterisks (permutation test,  $n = 19$ , tested-early:  $P < 0.001$ ; tested-late:  $P < 0.001$ ; difference:  $P = 0.013$ ). **(b)** Left and middle: cross-temporal decoding matrices of the early- (left) and late-tested (middle) items derived from training and testing on all timepoint combinations. Right: the difference between the decoding of the early- and late-tested item. Gray outline, timepoints of significantly lower decoding relative to both equivalent timepoints along the diagonal, which is taken as evidence for dynamic coding (permutation test,  $n = 19$ , cluster-defining threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Black outline (right), significantly higher decodability of the early-tested compared to the late-tested item (permutation test,  $n = 19$ , cluster-defining threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ).



**Figure 6** Attended and unattended WM items in early and late epochs and their relationship with behavioral performance. **(a)** Item decoding of the early- (blue) and late-tested items (red) during the first impulse epoch. Colored bars on top indicate significant decoding clusters of the corresponding items (permutation test,  $n = 19$ , cluster-defining threshold  $P < 0.05$ , corrected significance level  $P < 0.05$ ). Error shading, 95% CI of the mean. Boxplots and superimposed circles with error bars (mean and 95% CI of the mean) represent average decodability from 100 ms after stimulus onset until the end of the epoch. Significant average decoding and average differences between the decodability of the early and late item are marked by an asterisk (permutation test,  $n = 19$ , tested-early:  $P < 0.001$ ; tested-late:  $P = 0.045$ ; difference:  $P = 0.031$ ). Black and blue bar on the x axis shows onset of the first impulse. **(b)** Item decoding during the second impulse epoch; same conventions as **a**. Significant average decoding of the late item is marked by an asterisk (permutation test,  $n = 19$ ,  $P = 0.016$ ). Black and red bar on the x axis shows onset of the second impulse. **(c)** Left: boxplot, superimposed circles and error bars represent the differences in overall WM task performance between high- and low-decoding trials of early-tested items during the first impulse. Right: proportion of clockwise responses for high- and low-decoding trials as a function of the angular difference between the memory item and the probe. Error bars are 95% CI of the mean. Inset: boxplot and error bars for the difference in the slope parameter (a measure of memory precision) between high- and low-decoding trials. **(d)** As in **c** but for decoding the late-tested item during the late impulse. Significant differences in accuracy and/or precision between high- and low-decoding trials are highlighted by asterisks (permutation test,  $n = 19$ , early-tested item during first impulse:  $P = 0.038$  for accuracy; late-tested item during second impulse:  $P = 0.0404$  and  $P = 0.006$  for accuracy and precision, respectively, two-sided and one-sided for accuracy and precision tests, respectively). In boxplots, horizontal line indicates the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5x the interquartile range. Extreme values are shown separately (crosses).

### Experiment 3

We developed the impulse perturbation approach to reveal otherwise hidden neural states, without necessarily transforming the mnemonic representation<sup>4,21</sup>. This contrasts with other studies using retro-cues<sup>10,11,31</sup> or transcranial magnetic stimulation<sup>35</sup> to reactivate a latent item in working memory. However, to test whether our impulse stimulus actually did result in a behaviorally relevant transformation of the memory item (i.e., from a functionally latent to active state), we conducted an additional behavioral experiment ( $n = 20$  subjects). Adapting the design of Experiment 1, we now varied the presentation of the stimulus-onset asynchrony (SOA) between impulse and probe onset in Experiment 3 (SOA from 0 to 500 ms; **Supplementary Fig. 4a**). If the increase in impulse-specific decodability observed in both EEG experiments reflected a functional reactivation of an otherwise latent memory item, there should be a corresponding benefit to behavior.

A repeated-measures ANOVA provided no evidence for an effect of SOA ( $F_{4,76} = 1.184$ ,  $P = 0.325$ ). Uncorrected paired comparisons between the no-impulse condition (SOA 0 ms) and all other SOAs also provided no evidence for an impulse-specific effect on accuracy for any SOA (permutation test,  $n = 20$ , all  $P > 0.12$ ; **Supplementary Fig. 4b**). This suggests that our impulse stimulus was effective for ping-pong activity silent neural states without resulting in any behaviorally relevant transformation of the mnemonic representation.

### DISCUSSION

Recent theoretical models of WM predict a key role for activity-silent neural states in maintaining item-specific information<sup>4,17,18</sup>. This raises a particular challenge for contemporary neuroscience, which is dominated by measurement and analysis of neural activation states. Here we addressed this challenge using a perturbation approach to

reveal hidden neural states that encode the contents of WM. We showed that the response to an impulse stimulus faithfully reflected item-specific information in WM. We further demonstrated that the impulse response reflected both attended and unattended items in WM, yet recently forgotten information left no detectable traces in the hidden state. Behavioral modeling further suggested that the hidden-state coding determined the quality of information in WM.

Previous evidence from nonhuman primates shows that a neutral visual stimulus presented during the WM delay period can elicit distinct patterns of neural activity that depend on recent visual input<sup>36</sup>. Although the previous work could not deconfound previous sensory stimulation and WM proper, the observed effect helped motivate a dynamic coding model for WM<sup>4</sup>. According to this framework, distinct memoranda are associated with distinct changes in the neural response profile, which would be readable to downstream systems from the state-dependent response to a retrieval probe<sup>4,18</sup>. Crucially, WM depends on the maintenance of the item-specific neural response profile, rather than an explicit representation of an item in a persistent activity state. We now provide direct evidence for a WM-dependent impulse response decoupled from previous stimulation history and further demonstrate that this WM state is highly flexible and coupled to behavioral performance. The hidden state for a specific item can be rapidly cleared if it is no longer relevant to the task, providing a striking neural correlate of directed forgetting in WM.

Recent retro-cuing evidence suggests that prioritizing one WM item relative to other task-relevant items improves neural decoding of the cued item, whereas decoding of unattended items drops to chance levels even though the unattended information is still ultimately task-relevant and retrievable at the end of the trial<sup>10</sup>. Item-specific delay activity therefore seems to reflect the focus of attention, rather than WM per se<sup>31</sup>. The impulse response reported here clearly differed from the typical profile observed for decoding delay activity patterns. In Experiment 2, both attended and unattended items could be decoded from the impulse response of the hidden state as long as they were both still ultimately required for task performance. This suggests that if the information was successfully maintained in WM, there was a corresponding trace in the hidden state, irrespective of attentional priority. These results highlight the flexibility of WM, independently of attention-switching between specific items in WM. Activity states appear to track the focus of attention<sup>10,11,31</sup>, whereas hidden states, as revealed by the impulse response, more closely track the actual contents of WM.

Exactly how the proposed hidden state can be used for WM-guided behavior remains an important open question. Computationally, supervised learning could determine the mapping between the memory-dependent probe response and the correct behavioral response<sup>37</sup>, but such a learning strategy seems implausible for real-world behavior. Trial-and-error learning of arbitrary patterns does not seem a realistic model for WM, at least for humans. Instead, the inherent dynamics could establish a history-dependent match filter<sup>20</sup>, which would be capable of transforming probe input to a common decision signal (i.e., match versus no-match or, in our case, clockwise versus counterclockwise). In Myers *et al.*<sup>27</sup>, such a mechanism was shown to generate two distinct decision-related signals in an orientation detection task: a signed (i.e., directional) and unsigned difference signal, even though the signed difference was actually irrelevant to behavior in that task. A similar process could underpin WM encoding in hidden states. The hidden state could establish a flexible, task dependent circuit for WM-dependent decision-making<sup>38</sup>. When the probe stimulus is presented, the hidden state transforms the input to decision-relevant output: for example, direction of angular rotation.

However, because the impulse stimulus used in these experiments does not contain decision-relevant features, the impulse response reflects an input–output transformation of the arbitrary input.

It may be noted that, although the response to an arbitrary input is sufficient to read out the hidden state, it is unlikely to constitute an explicit reactivation of the memory representation. In contrast, retro-cueing can convert an unattended item to a prioritized state in preparation for recall<sup>22</sup>. Similarly, a recent transcranial magnetic stimulation study suggests that stimulation of the visual cortex can also render an item active from its latent state<sup>35</sup>. We find no evidence that our impulse stimulus reactivated the same pattern associated with stimulus processing. Moreover, a further behavioral experiment designed to test the possible behavioral consequences of our impulse stimulus provided no evidence that it interacted with the mnemonic representation. Rather, we argue that the impulse response simply ‘echoed’ the representational structure of the hidden state but did not drive an explicit transformation of latent memories to a prioritized state.

It has long been assumed that WM maintenance depends on persistent neural activity<sup>2</sup>. Instead, we propose that activity-silent neural states are sufficient to bridge memory delays. Activity-dependent transformations in hidden states determine the temporary coding properties of memory networks, i.e., dynamic coding<sup>4,36</sup>. WM decisions are made by the state-dependent response to subsequent input. However, WM is also classically associated with active manipulation of content in short-term memory<sup>1</sup>. We argue that such transformations are activity-dependent but that the results of the transformation can be maintained in short term memory via latent network states. This alternative account does not ignore previous evidence for decodable activity during mnemonic delays but rather attributes such evidence to focused attention<sup>35</sup>, periodic<sup>18</sup> or stochastic<sup>17</sup> updating, and/or response preparation<sup>8</sup>. Notably, our current results also showed that cue-directed forgetting can rapidly wipe the mnemonic representation from the hidden state. Rapid construction and dissolution of hidden states places important constraints on the basic mechanisms of hidden-state coding.

Although the present study addressed a specific model of WM, it is worth noting that the general impulse response approach for inferring otherwise silent neural states could also be particularly fruitful for exploring other tonic cognitive states, such as task set, attention and expectation. It is becoming increasingly apparent that we need to look beyond simple measures of neural activity and consider a richer diversity of neural states that underpin context-dependent behavior. Here we focus on perturbation to illuminate hidden states, but future work will also profit from more direct measures of functionally relevant hidden states (for example, synaptic efficacy, membrane potentials, extracellular transmitter concentrations). This will require more sophisticated measurements in awake behaving animals, coupled with noninvasive approaches like those described here for human studies.

## METHODS

Methods, including statements of data availability and any associated accession codes and references, are available in the [online version of the paper](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

We thank E. Spaak, A. Cravo and N. Myers for comments and advice and all our volunteers for their participation. We also thank the Biotechnology & Biological Sciences Research Council (BB/M010732/1 to M.G.S.) and the National Institute for Health Research Oxford Biomedical Research Centre Programme based at the Oxford University Hospitals Trust, Oxford University. The views expressed



are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health.

#### AUTHOR CONTRIBUTIONS

M.J.W., M.G.S., and E.G.A. designed the study. M.J.W. and J.J. collected the data. M.J.W. analyzed the data. M.J.W., M.G.S., E.G.A., and J.J. wrote the manuscript.

#### COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

- Baddeley, A. Working memory: looking back and looking forward. *Nat. Rev. Neurosci.* **4**, 829–839 (2003).
- Curtis, C.E. & D'Esposito, M. Persistent activity in the prefrontal cortex during working memory. *Trends Cogn. Sci.* **7**, 415–423 (2003).
- Goldman-Rakic, P.S. Cellular basis of working memory. *Neuron* **14**, 477–485 (1995).
- Stokes, M.G. 'Activity-silent' working memory in prefrontal cortex: a dynamic coding framework. *Trends Cogn. Sci.* **19**, 394–405 (2015).
- Watanabe, K. & Funahashi, S. Neural mechanisms of dual-task interference and cognitive capacity limitation in the prefrontal cortex. *Nat. Neurosci.* **17**, 601–611 (2014).
- Watanabe, K. & Funahashi, S. Prefrontal delay-period activity reflects the decision process of a saccade direction during a free-choice ODR task. *Cereb. Cortex* **17** Suppl 1: i88–i100 (2007).
- Miller, E.K., Erickson, C.A. & Desimone, R. Neural mechanisms of visual working memory in prefrontal cortex of the macaque. *J. Neurosci.* **16**, 5154–5167 (1996).
- Barak, O., Tsodyks, M. & Romo, R. Neuronal population coding of parametric working memory. *J. Neurosci.* **30**, 9424–9430 (2010).
- LaRocque, J.J., Lewis-Peacock, J.A., Drysdale, A.T., Oberauer, K. & Postle, B.R. Decoding attended information in short-term memory: an EEG study. *J. Cogn. Neurosci.* **25**, 127–142 (2013).
- Lewis-Peacock, J.A., Drysdale, A.T., Oberauer, K. & Postle, B.R. Neural evidence for a distinction between short-term memory and the focus of attention. *J. Cogn. Neurosci.* **24**, 61–79 (2012).
- Sprague, T.C., Ester, E.F. & Serences, J.T. Restoring latent visual working memory representations in human cortex. *Neuron* **91**, 694–707 (2016).
- Sreenivasan, K.K., Curtis, C.E. & D'Esposito, M. Revisiting the role of persistent neural activity during working memory. *Trends Cogn. Sci.* **18**, 82–89 (2014).
- Buonomano, D.V. & Maass, W. State-dependent computations: spatiotemporal processing in cortical networks. *Nat. Rev. Neurosci.* **10**, 113–125 (2009).
- Barak, O. & Tsodyks, M. Working models of working memory. *Curr. Opin. Neurobiol.* **25**, 20–24 (2014).
- Murray, J.D. *et al.* Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proc. Natl. Acad. Sci. USA* **114**, 394–399 (2017).
- Fujisawa, S., Amarasingham, A., Harrison, M.T. & Buzsáki, G. Behavior-dependent short-term assembly dynamics in the medial prefrontal cortex. *Nat. Neurosci.* **11**, 823–833 (2008).
- Lundqvist, M. *et al.* Gamma and beta bursts underlie working memory. *Neuron* **90**, 152–164 (2016).
- Mongillo, G., Barak, O. & Tsodyks, M. Synaptic theory of working memory. *Science* **319**, 1543–1546 (2008).
- Hempel, C.M., Hartman, K.H., Wang, X.-J., Turrigiano, G.G. & Nelson, S.B. Multiple forms of short-term plasticity at excitatory synapses in rat medial prefrontal cortex. *J. Neurophysiol.* **83**, 3031–3041 (2000).
- Sugase-Miyamoto, Y., Liu, Z., Wiener, M.C., Optican, L.M. & Richmond, B.J. Short-term memory trace in rapidly adapting synapses of inferior temporal cortex. *PLoS Comput. Biol.* **4**, e1000073 (2008).
- Wolff, M.J., Ding, J., Myers, N.E. & Stokes, M.G. Revealing hidden states in visual working memory using electroencephalography. *Front. Syst. Neurosci.* **9**, 123 (2015).
- Griffin, I.C. & Nobre, A.C. Orienting attention to locations in internal representations. *J. Cogn. Neurosci.* **15**, 1176–1194 (2003).
- Landman, R., Spekreijse, H. & Lamme, V.A.F. Large capacity storage of integrated objects before change blindness. *Vision Res.* **43**, 149–164 (2003).
- Harrison, S.A. & Tong, F. Decoding reveals the contents of visual working memory in early visual areas. *Nature* **458**, 632–635 (2009).
- Worden, M.S., Foxe, J.J., Wang, N. & Simpson, G.V. Anticipatory biasing of visuospatial attention indexed by retinotopically specific alpha-band electroencephalography increases over occipital cortex. *J. Neurosci.* **20**, RC63 (2000).
- Saproo, S. & Serences, J.T. Spatial attention improves the quality of population codes in human visual cortex. *J. Neurophysiol.* **104**, 885–895 (2010).
- Myers, N.E. *et al.* Testing sensory evidence against mnemonic templates. *eLife* **4**, e09000 (2015).
- Zhang, W. & Luck, S.J. Discrete fixed-resolution representations in visual working memory. *Nature* **453**, 233–235 (2008).
- Bays, P.M. & Husain, M. Dynamic shifts of limited working memory resources in human vision. *Science* **321**, 851–854 (2008).
- Murray, A.M., Nobre, A.C. & Stokes, M.G. Markers of preparatory attention predict visual short-term memory performance. *Neuropsychologia* **49**, 1458–1465 (2011).
- LaRocque, J.J., Lewis-Peacock, J.A. & Postle, B.R. Multiple neural states of representation in short-term memory? It's a matter of attention. *Front. Hum. Neurosci.* **8**, 5 (2014).
- Olivers, C.N.L., Peters, J., Houtkamp, R. & Roelfsema, P.R. Different states in visual working memory: when it guides attention and when it does not. *Trends Cogn. Sci.* **15**, 327–334 (2011).
- Souza, A.S. & Oberauer, K. In search of the focus of attention in working memory: 13 years of the retro-cue effect. *Atten. Percept. Psychophys.* **78**, 1839–1860 (2016).
- van Ede, F., Niklaus, M. & Nobre, A.C. Temporal expectations guide dynamic prioritization in visual working memory through attenuated  $\alpha$  oscillations. *J. Neurosci.* **37**, 437–445 (2017).
- Rose, N.S. *et al.* Reactivation of latent working memories with transcranial magnetic stimulation. *Science* **354**, 1136–1139 (2016).
- Stokes, M.G. *et al.* Dynamic coding for cognitive control in prefrontal cortex. *Neuron* **78**, 364–375 (2013).
- Mante, V., Sussillo, D., Shenoy, K.V. & Newsome, W.T. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- Martínez-García, M., Rolls, E.T., Deco, G. & Romo, R. Neural and computational mechanisms of postponed decisions. *Proc. Natl. Acad. Sci. USA* **108**, 11626–11631 (2011).

## ONLINE METHODS

**Participants.** Thirty healthy adults (13 female, mean age 24.9 years, range 18–38 years) were included in the analyses of Experiment 1, 19 adults (10 female, mean age 24.7 years, range 18–39 years) in Experiment 2 and 20 adults in Experiment 3 (13 female, mean age 21, range 18–29 years). During data collection and preprocessing, four additional participants in Experiment 1, one additional participant of Experiment 2 and six additional participants of Experiment 3 were excluded from all analyses due to either low average performance on the memory task (below 60% accuracy) or excessive eye movements (more than 30% of trials contaminated). No statistical methods were used to predetermine sample sizes, but our sample sizes are similar to those reported in previous publications<sup>21,28</sup>. All participants of Experiment 1 and 2 received monetary compensation of £10/h, and participation in Experiment 3 contributed to course credits. All participants gave written informed consent. Experiments 1 and 2 were approved by the Central University Research Ethics Committee of the University of Oxford and Experiment 3 was approved by the Departmental Ethical Committee of the University of Groningen.

**Apparatus and stimuli.** The experimental stimuli were generated and controlled by Psychtoolbox<sup>39</sup>, a freely available Matlab extension. The stimuli were presented on a 23-inch (58.42-cm) screen running at 100 Hz and a resolution of 1,920 by 1,080 in Experiment 1; on a 22-inch (55.88-cm) screen at a resolution of 1,680 by 1,050 in Experiment 2; and on a 19-inch (48.26-cm) CRT screen running at 100 Hz and a resolution of 1,280 by 1,024 in Experiment 3. Viewing distance was set at 64 cm in Experiment 1, 67.5 cm in Experiment 2 and approximately 60 cm (not controlled) in Experiment 3, to ensure that the visual angles of stimuli were the same across experiments even though the screen parameters were different. A standard keyboard was used for response input by the participants.

All reported stimuli were the same in all experiments, unless explicitly mentioned otherwise. A gray background (RGB = 128, 128, 128; 20.5 cd/m<sup>2</sup>; 28.6 cd/m<sup>2</sup> in Experiment 3) was maintained throughout the experiments. A black fixation dot with a white outline (0.242°) was presented in the center of the screen throughout all trials. Memory items and memory probes were sine-wave gratings presented at 20% contrast, with a diameter of 6.69° and spatial frequency of 0.65 cycles per degree. The phase was randomized within and across trials. The memory items were presented at 6.69° eccentricity, and for each trial the orientations were randomly selected without replacement from a uniform distribution of orientations. The impulse stimulus was three adjacent 'bull's-eyes' in Experiment 1. Each bull's-eye was of the same size and spatial frequency as the memory items. To reduce strain on the eyes, and to minimize forward masking in Experiment 3, the impulse stimulus in Experiments 2 and 3 consisted of three adjacent white circles. In Experiment 1 and 2 the probes had the same contrast and spatial frequency as the memory items and were presented in the center of the screen. In Experiment 3 the probe screen included a high contrast black and white square-wave grating in the center and two white lateralized circles on the outside (the same location and size as the preceding lateral impulse circles). The angle differences between a memory item and the corresponding memory probe were uniformly distributed across seven angle differences in Experiment 1 ( $\pm 3^\circ$ ,  $\pm 7^\circ$ ,  $\pm 12^\circ$ ,  $\pm 18^\circ$ ,  $\pm 25^\circ$ ,  $\pm 33^\circ$ ,  $\pm 42^\circ$ ), six angle differences in Experiment 2 ( $\pm 5^\circ$ ,  $\pm 10^\circ$ ,  $\pm 16^\circ$ ,  $\pm 24^\circ$ ,  $\pm 26^\circ$ ,  $\pm 32^\circ$ ,  $\pm 40^\circ$ ) and a single angle difference ( $\pm 16^\circ$ ) in Experiment 3.

**Procedure.** *Experiment 1.* Participants completed a retro-cue visual working memory task. Each trial began with the onset of a fixation dot at the center of the screen. After 1,000 ms, the memory item array was shown for 250 ms, consisting of two randomly oriented low-contrast gratings left and right of fixation. After a delay of 800 ms an arrow was shown for 200 ms in the center of the screen, pointing either to the left or to the right, and thus cueing which of the two previously presented items would be tested. The number of left and right cued trials was equal and the order was randomized for each participant. The impulse stimulus was presented for 100 ms, 900 ms after the offset of the retro-cue. After another delay of 400 ms, the memory probe was shown for 250 ms. Participants were instructed to indicate if the orientation of the probe relative to the orientation of the memory item was rotated clockwise by pressing the 'm' key with the right index finger or counter-clockwise by pressing the 'c' key with the left index finger. A high or low frequency feedback tone was played after response, indicating if the answer was correct or incorrect, respectively. The next trial started

within 400–700 ms (determined randomly). Participants completed 1,344 trials in total, which took approximately 3 h (including breaks). Trial conditions were randomized across the whole session. See **Figure 1a** for a trial schematic.

*Experiment 2.* Participants completed a visual working memory task in which two items were serially tested. The experiment began by instructing the participant which of the two memory items would be tested early and which one would be tested late. This rule never changed within a session. Each trial began with the onset of a fixation dot at the center of the screen. After 1,000 ms, the memory item array was shown for 250 ms, consisting of two randomly oriented low-contrast gratings left and right of fixation. After a delay of 950 ms, the first impulse was presented for 100 ms. After a delay of 500 ms, the first memory probe was presented for 250 ms, probing the first item. The response input was the same as in Experiment 1. After a fixed delay of 1,750 ms after the offset of the first probe, the second impulse was shown for 100 ms. Following a delay of 400 ms, the second memory probe was presented for 250 ms, probing the late-tested item. After the second response, two feedback tones were played, one for each response, separately indicating whether the first and second answers were correct. Participants completed two sessions of the task on two separate days, separated by approximately 1–2 weeks. The testing order of the memory items was fixed within each session and switched between sessions (i.e., if the left item tested first in one session, the right item was tested first in the other session). The order of the testing rule between sessions (i.e., whether the left item was tested first in the first session or in the second session) was counterbalanced across participants (odd-numbered participants were tested on the left first, even-numbered were tested on the right first). Each session consisted of 864 trials and lasted approximately 3 h including breaks. See **Figure 3a** for a trial schematic.

*Experiment 3.* The task was almost the same as in Experiment 1, including the same timings of the memory items, cue, probe and overall trial duration. The one key difference was the timing of the impulse stimulus. While the delay between cue offset and probe onset was held constant at 1,400 ms across all trials (the same as in Experiment 1), the SOA between impulse and probe onset was 0, 50, 100, 250 or 500 ms (determined pseudorandomly across the session). No impulse was shown in the 0 ms SOA condition. The impulse remained on the screen until the probe stimulus was presented. This was to ensure the least possible interference from the impulse on probe processing (i.e., rapid onset and offset of the white circles immediately before probe presentation could deteriorate probe visibility), as well as keeping the different SOA conditions as similar as possible (longer SOA would include an additional offset). Participants completed 280 trials (approximately 30 min). See **Supplementary Figure 4a** for a trial schematic.

Data collection and analyses were not performed blind to the conditions of the experiments. Due to the within-subject design in all three experiments, randomization of conditions between subjects was not applicable.

**EEG acquisition.** The EEG signal was acquired from 61 Ag/AgCl sintered electrodes (EasyCap, Herrsching, Germany) laid out according to the extended international 10–20 system. Data was recorded at 1,000 Hz using a NeuroScan SynAmps RT amplifier and Scan 4.5 software in Experiment 1 and Curry 7 software in Experiment 2 (Compumedics NeuroScan, Charlotte, NC). The anterior midline frontal electrodes (AFz) served as the ground. Bipolar electrooculography (EOG) was recorded from electrodes placed above and below the right eye, to the left of the left eye and to the right of the right eye. The impedances of all electrodes were kept below 5 k. Online, the EEG was referenced to the right mastoid and filtered using a 200-Hz low-pass filter.

**EEG preprocessing.** Offline, the data was re-referenced to the average of both mastoids, down-sampled to 500 Hz and bandpass filtered (0.1 Hz high-pass and 40 Hz low-pass) using EEGLAB<sup>40</sup>. The data was then epoched to the onset of the memory items and the impulse. In Experiment 1, the memory item epoch was from –200 ms to 1,050 ms, relative to onset, and in Experiment 2 from –200 ms to 1,200 ms. The impulse epochs were from –200 ms to 500 ms relative to onset in both experiments. Additionally, for the purpose of artifact rejection, including rejection of trials containing saccadic eye movements before the time of interest (see below), the cue segment in Experiment 1 was also epoched (–200 ms to 1,100 ms).

Subsequent artifact detection and trial rejection focused exclusively on the 17 posterior channels that were included in the analyses (P7, P5, P3, P1, Pz, P4, P6,

P8, PO7, PO3, POz, PO4, PO8, O1, Oz and O2) and the EOGs. Each trial of each epoch was individually visually inspected for blinks, saccades and nonstereotyped artifacts. Trials from individual epochs were rejected from analyses involving that epoch if it contained any of the above-mentioned artifacts. Furthermore, impulse-epoch trials were also excluded from corresponding analyses if the EOG signal suggested that saccades occurred during any of the previous epochs of that trial. In Experiment 1, this exclusion procedure was applied to the cue-epoch as well. In Experiment 2, late-impulse trials were also excluded if no response was registered for the preceding probe. For the decoding analyses, each epoch was baselined using the average signal from -200 ms to 0 ms before stimulus onset. The multivariate data were also demeaned at each timepoint by subtracting the average voltage for all posterior channels included in the analyses.

**Time-frequency decomposition and lateralization analysis.** In order to explore alpha power (8–12-Hz) lateralization<sup>25,41</sup>, the spectral power from 6 to 16 Hz (in steps of 0.5 Hz) of the EEG signal was computed using Hanning tapers with time-windows of 5 cycles per frequency (in steps of 10 ms) using the Matlab toolbox FieldTrip<sup>42</sup>. We included the whole experimental trial, ranging from 1,000 ms before memory item onset until 1,500 ms after (second) probe onset (-1,000 to +4,150 ms relative to memory items in Experiment 1 and -1,000 to +5,800 ms relative to memory items in Experiment 2). The power was log-transformed, and lateralization was computed by subtracting the average power of the ipsilateral posterior electrodes from the average power of the contralateral posterior electrodes in relation to the cued memory item in Experiment 1 and to the early-tested item in Experiment 2 (P7, P5, P3, P1, PO7, PO3 and O1 versus P8, P5, P6, P4, P2, PO8, PO4 and O2).

Significant clusters of lateralization were determined using a cluster-corrected nonparametric sign-permutation test<sup>43</sup>. In both experiments, the whole trial was included in this analysis (-100 to +3,150 ms relative to memory items onset in Experiment 1 and -100 to +4,800 ms in Experiment 2).

**Orientation decoding.** To test whether the activity pattern of the posterior EEG channels of interest contained orientation-specific activity, we used the Mahalanobis distance<sup>44</sup> to compute the trial-wise distances between the full range of possible orientations and quantify to what extent the computed distances adhered to the parametric circular space of the orientations<sup>11</sup>. This approach is an extension of the pairwise distance approach we used before<sup>21</sup> and is conceptually similar to the population tuning curve model<sup>26</sup>.

The left and right presented items were decoded separately and independently within each participant and experimental session. All 17 posterior channels (see above) were used for all decoding analyses. The procedure followed a leave-one-trial-out cross-validation approach to compute the trial-wise decodability of the orientation of interest. The activity pattern of a single test trial at a particular timepoint was compared to the pattern of all other trials at the same timepoint. These were averaged into 12 orientation bins relative to the orientation of the test trial, each containing trials with orientations within a range of 30° and centered around -75°, -60°, -45°, -30°, -15°, 0°, 15°, 30°, 45°, 60°, 75° and 90°. The Mahalanobis distances between the test trial and each orientation bin was computed using the covariance estimated from all trials, excluding the test trial, using a shrinkage estimator<sup>45</sup>. To simplify visualization and interpretation, the 12 resulting distances were mean-centered and the sign was reversed, resulting in a visual representation of a tuning curve. Higher values correspond to greater relative similarity between the test trial and the averaged train trials within a particular orientation bin, and lower values correspond to greater dissimilarity.

Next, the vector means of the tuning curves were computed<sup>11</sup>. First, the cosine of the center of each orientation bin ( $\theta$ ) was rescaled to the range -180 to 180. It was then multiplied by the corresponding sign-reversed distances ( $d(\theta)$ ) before the mean of the resulting 12 values was taken, which made up the decoding accuracy ( $da$ ).

$$da = \text{mean}(d(\theta) \cos(2\theta)) \quad (1)$$

A high value reflects evidence for orientation tuning: the difference between the test trial and train trials with a similar orientation is smaller than between the test trial and train trials with different orientations. This procedure was repeated for all trials and all timepoints. See **Supplementary Software** for the custom Matlab function used to decode orientations using Mahalanobis distance.

The decoding values were averaged over all trials and smoothed over time with a Gaussian smoothing kernel (s.d. = 16 ms) for visualization and time-resolved significance testing.

Cluster-corrected sign-permutation significance tests were carried out within the memory items epoch (0 to 1,050 ms in Experiment 1; 0 to 1,200 ms in Experiment 2) and impulse epochs separately (0 to 500 ms in both experiments), to explore the significant decoding time-course. Additionally, to assess the overall decodability within an epoch, the decoding values were averaged over time (from 100 ms after stimulus onset until the end of the epoch) and then submitted to a two-sided permutation test.

**Relationship between behavior and decoding.** The trial-wise average decoding scores after memory items presentation (100 to 1,050 ms in Experiment 1 and 100 to 1,200 ms in Experiment 2) and impulse presentation (100 ms to 500 ms) was median split. Nonresponse trials (to the early probe in Experiment 2) were excluded from this analysis. The average behavioral accuracies of high- and low-decoding trials were statistically compared using a two-sided permutation test.

**Behavioral modeling.** To further explore the relationship between WM task performance and trial-wise decoding, we modeled the behavioral performance as a function of the difference in degrees between the orientation of the memory item and the probe using the following model, which was fit to each participant separately<sup>30</sup>.

$$y = \lambda + \frac{(1 - 2\lambda)}{2} \times \text{erfc}\left(\frac{-\beta}{\sqrt{2}}(x - \alpha)\right) \quad (2)$$

where  $\text{erfc}$  is the complementary Gaussian error function,  $\lambda$  is the asymptote,  $\beta$  is the slope and  $\alpha$  is the threshold/bias parameter. The model fitting was performed using the Palamedes Matlab toolbox (<http://www.palamedestoolbox.org/>). The asymptote represents the guess rate, where a higher value reflects a higher probability that no information about the probed item is maintained in WM, resulting in a higher probability for mistakes even when the angular difference between the probe and the memory item is large. The slope is interpreted as the memory precision, where a high precision reflects a relatively high proportion of correct responses at small degree rotations between the probe and memory item. The asymptote and slope parameters were both unconstrained across the high- and low-decoding conditions. A single bias parameter was used, which was included (instead of fixing it at 0) because cumulative-likelihood tests<sup>46</sup> showed better model fits for all cases (Experiment 1:  $n = 30$ ,  $\chi^2_{30} = 135.978$ ,  $P < 0.001$ ; Experiment 2,  $n = 19$ , early accuracy:  $\chi^2_{19} = 215.351$ ,  $P < 0.001$ ; late accuracy:  $\chi^2_{19} = 33.69$ ,  $P = 0.02$ ).

The unconstrained model parameters (slope and asymptote) were subsequently compared between high- and low-decoding trials. Since the behavioral modeling was carried out as a direct follow up to the average accuracy effects observed in both experiments (two-sided tests), we had clear expectations about the directionality of the effects. For the positive relationship between decoding and accuracy observed for the cued item in Experiment 1 and for both tests in Experiment 2, we expected that decoding should have a negative relationship with the guess rate (i.e., lower guess rates for higher decoding) and/or a positive relationship with precision (higher precision for higher decoding) and vice versa for the negative accuracy effect of the uncued item in Experiment 1. Therefore, all tests of model parameter comparisons between high- and low-decoding trials were one-sided.

**Cross-temporal decoding.** We also explored the cross-temporal dynamics of stimulus processing and maintenance as a function of item priority in Experiment 2 and the cross-generalization between impulse and memory presentation epochs in both experiments. The decoding approach was the same as described above, except classifiers trained at each time point were tested at every other time point, resulting in two-dimensional cross-temporal decoding matrices<sup>47</sup>.

If the decoding patterns are stationary, it should not matter whether training/testing is performed using the same time points. In contrast, decoding often appears dynamic: training and testing on the same timepoints results in higher decoding scores than training and testing on different timepoints (i.e., minimal cross-temporal generalization). We tested for this hallmark feature of dynamic coding using a nonparametric test used previously<sup>27</sup>. The decodability at each

cross-temporal timepoint  $t_{x,y}$  was compared to the pair of decodabilities at the corresponding within timepoints ( $t_{x,x}$  and  $t_{y,y}$ ) with two separate permutation tests. A significant difference in both was taken as evidence for dynamic coding. Timepoints of significant dynamic coding were corrected for multiple comparisons using a two-dimensional cluster-based permutation test.

**Significance testing.** To determine statistical significance, we used the nonparametric sign-permutation test<sup>43</sup> (with one exception; see discussion of ANOVA below), which does not make assumptions about the underlying distribution. Since the null hypotheses of all tests corresponded to no effect (i.e., no difference in power lateralization, no difference in decodability, etc.), the sign of the data of each participant was randomly flipped with a probability of 50% 50,000 times. The resulting distribution was used to derive the  $P$  value of the null hypothesis that the mean effect was equal to 0. All tests were two-sided, unless otherwise stated.

For time-series and frequency data, the above procedure was repeated for each timepoint and frequency (when applicable). To correct for multiple comparisons over time and/or frequencies, a cluster-based permutation test was subsequently used, with 50,000 permutations (5,000 for cross-temporal decoding, due to computer memory limitations) and using a cluster-forming threshold and cluster significance threshold of  $P < 0.05$ . Tests concerning the average of specific time-windows (including decoding-behavior relationships) were performed to test unique and independent hypotheses, and therefore no corrections were applied. The sample sizes for all tests were  $n = 30$  in Experiment 1,  $n = 19$  in Experiment 2 and  $n = 20$  in Experiment 3. The 95% confidence intervals of the error bars were determined by bootstrapping from the corresponding data 50,000 times.

The boxplots used in our figures follow the standard conventions. The middle line represents the median, the box the first and third quartile, and the whiskers all data within  $1.5\times$  the interquartile range of the lower and upper quartile. Where appropriate, data points outside this range are displayed individually (small crosses).

A repeated measures ANOVA was used to analyze the behavioral data of Experiment 3. The normality and equal-variances assumptions were tested with the Shapiro-Wilk test of normality and Mauchly's test of sphericity,

respectively. Neither test provided evidence for assumption violations of the data. A **Supplementary Methods Checklist** is available.

**Data availability.** The data that support the finding of this study are publically available at <http://datasharedrive.blogspot.co.uk/2017/03/dynamic-hidden-states-underlying.html>. All necessary task/condition information has been provided within a self-contained format, as specified in the *OECD Principles and Guidelines for Access to Research Data from Public Funding*<sup>48</sup>.

**Code availability.** The custom Matlab orientation decoding function is provided with the paper (**Supplementary Software**). All complete custom Matlab routines used to generate the figures of this paper are available at <http://datasharedrive.blogspot.co.uk/2017/03/dynamic-hidden-states-underlying.html>.

39. Brainard, D.H. The psychophysics toolbox. *Spat. Vis.* **10**, 433–436 (1997).
40. Delorme, A. & Makeig, S. EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J. Neurosci. Methods* **134**, 9–21 (2004).
41. Schneider, D., Mertes, C. & Wascher, E. The time course of visuo-spatial working memory updating revealed by a retro-cuing paradigm. *Sci. Rep.* **6**, 21442 (2016).
42. Oostenveld, R., Fries, P., Maris, E. & Schoffelen, J.M. FieldTrip: open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput. Intell. Neurosci.* **2011**, 156869 (2011).
43. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *J. Neurosci. Methods* **164**, 177–190 (2007).
44. De Maesschalck, R., Jouan-Rimbaud, D. & Massart, D.L. The Mahalanobis distance. *Chemometr. Intell. Lab. Syst.* **50**, 1–18 (2000).
45. Ledoit, O. & Wolf, M. Honey, I shrunk the sample covariance matrix. *J. Portfolio Management* **30**, 110–119 (2004).
46. Claessens, P.M.E. & Wagemans, J. A Bayesian framework for cue integration in multistable grouping: Proximity, collinearity, and orientation priors in zigzag lattices. *J. Vis.* **8**, 1–23 (2008).
47. King, J.-R. & Dehaene, S. Characterizing the dynamics of mental representations: the temporal generalization method. *Trends Cogn. Sci.* **18**, 203–210 (2014).
48. Pilat, D. & Fukasaku, Y. OECD principles and guidelines for access to research data from public funding. *Data Sci. J.* **6**, OD4–OD11 (2007).