

1 Dynamic hidden states underlying working memory guided behaviour

2

3 Michael J. Wolff^{1,2}, Janina Jochim¹, Elkan G. Akyürek² & Mark G. Stokes¹

4

5 ¹ Department of Experimental Psychology, University of Oxford, Oxford, United Kingdom

6 ² Department of Experimental Psychology, University of Groningen, Groningen, The
7 Netherlands

8

9

10 Corresponding author: Mark G. Stokes, mark.stokes@psy.ox.ac.uk

11

12

13 Summary of Main Finding: Wolff and colleagues show that ‘activity-silent’ brain states play an
14 important role in working memory. Using a novel perturbation method to ‘ping the brain’, they
15 uncover hidden neural states that reflect temporary information held in mind, and predict
16 memory performance. They argue that dynamic hidden states could underpin working memory.

17 Abstract

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

30

31 Working memory (WM) is a core cognitive function critical for flexible, intelligent behaviour¹.
32 Until recently, it was widely assumed that information is maintained in WM by maintaining
33 specific activity states that represent the specific memoranda^{2,3}. However, accumulating
34 evidence increasingly shows that successful maintenance in WM is not strictly dependent on an
35 unbroken chain of corresponding delay activity⁴, and that item-specific activity states could
36 reflect other cognitive processes. For example, in monkey studies persistent activity ramps up
37 with expectation of the probe⁵⁻⁸. Similarly, in the human it has been shown that unattended WM
38 content is not reflected in the neural signal, even when it is still clearly maintained⁹⁻¹¹. Evidence
39 for WM in the absence of persistent delay activity suggests that WM can be maintained in
40 ‘activity silent’ neural states⁴.

41 Recent theories acknowledge that brain activity is highly dynamic, even when the contents of
42 working memory remain stable¹². Multiple neurophysiological mechanisms could underlie such
43 dynamics¹³⁻¹⁵. According to a dynamic coding model of WM⁴, behaviourally relevant sensory
44 input drives a memory item-specific neural response, which triggers an item-specific change in
45 the functional state of the system. Depending on the precise neural mechanism, this functional
46 state could be activity-silent (e.g., short-term synaptic plasticity^{14,16-19}), and maintained
47 throughout the memory delay to serve as the neural context for subsequent processing. Items in
48 WM would be read-out via the context-dependent response to a probe stimulus during recall^{13,20}.
49 Crucially, this model predicts that dynamic hidden states are constructed when new information
50 is encoded, and dissolved as soon as it is forgotten. This model also predicts that dynamic hidden
51 states should determine the quality of a representation maintained in WM.

52 To probe hidden neural states, we developed a functional perturbation approach to ‘ping the
53 brain’. Analogous to the idea of active sonar (or echolocation), the response to a well-
54 characterised impulse stimulus can be used to infer the current state of the system^{4,13}. We
55 recently validated this general approach using non-invasive electroencephalography (EEG) in a
56 proof of principle study²¹. The presentation of a high contrast, neutral visual stimulus evoked
57 neural activity that clearly discriminated the previously presented visual stimulus. Here, we
58 exploit this approach to track the functional dynamics of hidden states for WM.

59 Across two experiments, we show that the content of WM can be decoded from the impulse
60 response during the maintenance interval, while forgotten information leaves effectively no
61 trace. In Experiment 2, we also demonstrate robust hidden-state representation for unattended
62 content in WM, providing a plausible mechanism for maintenance that is independent of the
63 activity associated with the focus of attention. Finally, we also find evidence that the quality of
64 working memory varies with the decodability of these hidden states.

Results

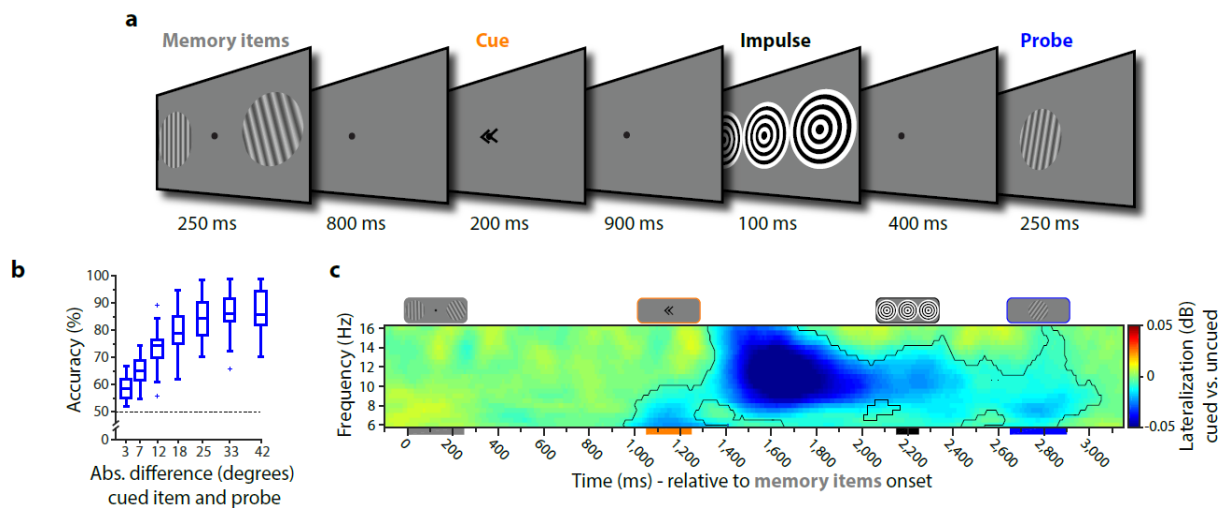
65

Experiment 1:

66

67 In Experiment 1, 30 human participants performed a visual WM task while EEG was recorded.
 68 At the beginning of each trial (see Fig. 1a), two memory items were presented, but a
 69 retrospective cue (retro-cue) presented during the delay instructed participants which item would
 70 actually be probed^{22,23}. The other item could be simply forgotten. The retro-cue in this design is
 71 essential to differentiate WM from basic stimulation history²⁴. During a subsequent memory
 72 delay, we then presented a high contrast “impulse” stimulus. Memory performance for the cued
 73 item was tested after the impulse by a centrally presented memory probe (Fig. 1b). Time-
 74 frequency decomposition of lateralised activity in posterior sensors (Fig. 1c) shows significant
 75 lateralization in the alpha range (8-12 Hz) after the presentation of the cue (permutation test, $n =$
 76 30, $p < 0.001$, corrected, cluster-forming threshold $p < 0.05$). This pattern is consistent with a
 77 shift in spatial attention²⁵ according to the retro-cue, which confirms that the cue manipulation
 78 was effective.

79



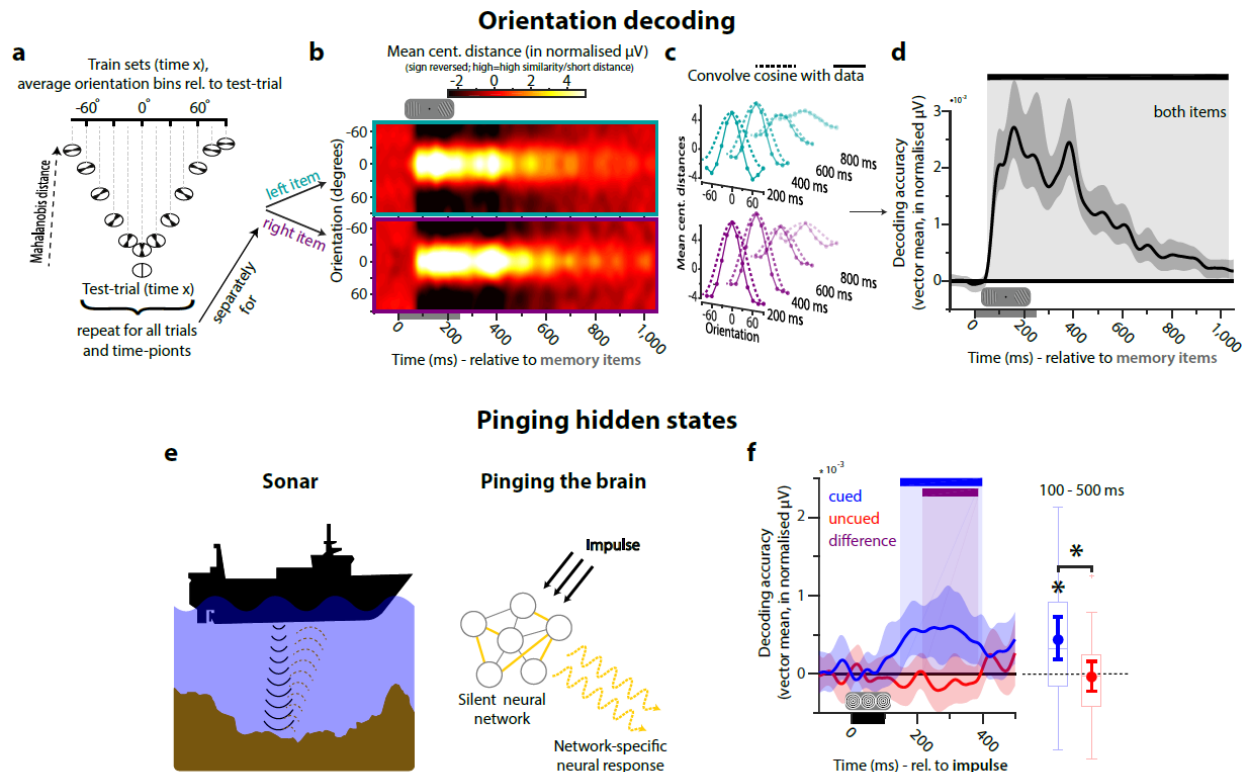
80

81

82 **Figure 1. Experiment 1 task structure, behavioural performance and attention-related**
 83 **alpha band activity. a.** Trial schematic. Two memory items were presented (randomly oriented
 84 grating stimuli), and participants were instructed to memorize both orientations. A retro-cue then
 85 indicated which item would actually be tested at the end of the current trial (100% valid). The

86 impulse stimulus (high contrast, task-irrelevant visual input) was then presented during the
87 subsequent delay while participants should have only the cued item in WM. At the end of the
88 trial, a forced-choice probe was presented at the centre of the screen. Participants indicated
89 whether the probe was rotated clockwise or anti-clockwise relative to the orientation of the cued
90 item. **b.** Boxplots show WM accuracy as a function of the absolute angular difference (in
91 degrees) between the memory item and the probe. Data points outside of the 1.5 * interquartile
92 range are shown separately (small crosses). **c.** Time-frequency representation of the difference
93 between the contra- and ipsilateral posterior electrodes relative to the cued hemifield. The
94 highlighted cluster in the alpha frequency band (8-12 Hz) indicates significant contralateral
95 desynchronization (permutation test, $n = 30$, cluster-forming threshold $p < 0.05$, corrected
96 significance level $p < 0.05$). The coloured bars under the x-axis represent the timings of the
97 corresponding stimuli illustrated on top.

98 Decoding parametric memory items: To decode the memory items used in this experiment, we
99 developed a parametric variant of distance-based discrimination (see Online Methods, Fig. 2a-d).
100 As shown in Fig 2a, this capitalises on the parametric structure of the stimulus space ²⁶, whilst
101 maintaining the statistical advantages of the Mahalanobis distance metric used in previous
102 EEG/MEG decoding studies ^{21,27} (see Online Methods). To summarise briefly here: for a given
103 trial, we compare the activity pattern across electrodes to the corresponding activity pattern
104 observed in the remaining trials, averaged by orientation-difference to the test trial (at a bin
105 width of 30 degrees). This procedure is repeated for all trials and all time-points. If the pattern of
106 activity contains information about item orientation, we expect greater pattern dissimilarity (i.e.,
107 Mahalanobis distance) at larger angular differences. Fig. 2b shows distance as a function of
108 reference angle and time after the presentation of the left and right item separately (upper/lower
109 respectively). Distance values were then converted into a decoding accuracy score (Fig. 2c) and
110 averaged across both items at each time-point (Fig. 2d). Item orientation could be decoded from
111 56 ms until 1026 ms after onset (permutation test, $n = 30$, $p < 0.001$ (corrected), cluster-forming
112 threshold $p < 0.05$). This is consistent with previous empirical evidence that EEG is sufficiently
113 sensitive to detect subtle differences in scalp-level activity patterns associated with different
114 stimulus orientation ²¹. The current decoding results further validate the utility of multivariate
115 pattern analysis for two simultaneously presented orientation gratings. For completeness, we also
116 decode item-specific orientation during the retro-cue epoch (Supplementary Fig. 1).



117

118 **Figure 2. Orientation decoding in EEG and pinging hidden states of WM. a-d.** Decoding
 119 procedure. **a.** The dissimilarity in the neural pattern between a single trial and all other trials is
 120 computed as a function of orientation difference (binned: 30 degrees). **b.** Average distance to
 121 template of all trials for each time-point during and after memory item presentation, plotted
 122 separately for the left and the right memory item (upper/lower respectively). Distances are mean
 123 centred and sign reversed (high = small distance/high similarity) for visualization. **c.** A cosine is
 124 convolved with the data. **d.** The vector mean of the convolved tuning curves (i.e., decoding
 125 accuracy) over time, averaged over left and right items. The black bar indicates significant
 126 decoding (permutation test, $n = 30$, cluster-forming threshold $p < 0.05$, corrected significance
 127 level $p < 0.05$). Error shading is the 95 % C.I. of the mean. **e.** Pinging hidden states. Analogy to
 128 active sonar: differences in hidden state are inferred from differences in the measured response to
 129 a well-characterised impulse. **f.** Decoding results in the impulse epoch. The blue bar indicates
 130 significant decoding of the cued item. The purple bar indicates significant difference in
 131 decodability between the cued and uncued item (permutation test, $n = 30$, cluster-forming
 132 threshold $p < 0.05$, corrected significance level $p < 0.05$). Error shading is the 95 % C.I. of the
 133 mean. The boxplots and superimposed circles with error-bars (mean and 95 % C.I. of the mean)

134 represent average decoding from 100 to 500 ms after impulse onset. Data points outside of the
135 1.5 * interquartile range are shown separately (small crosses). Significant average decoding and
136 significant difference in average decodability between the cued and uncued item are marked by
137 asterisks (permutation test, $n = 30$, $p < 0.05$).

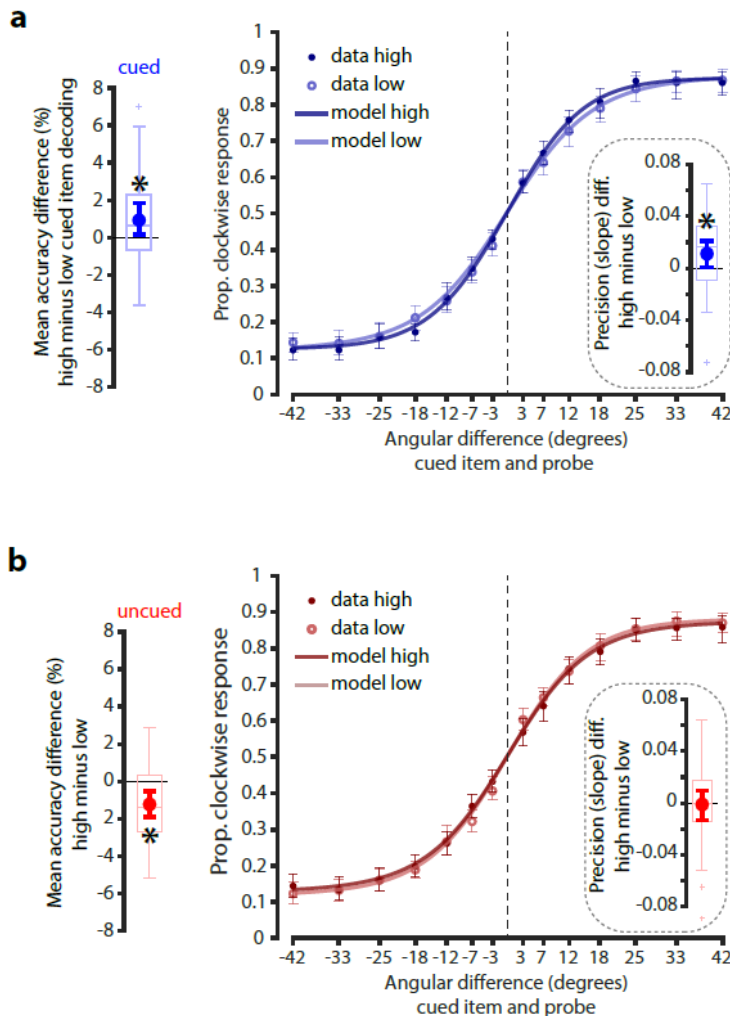
138 Pinging hidden states: According to the dynamic coding framework, we hypothesised that the
139 input/output mapping of neural circuits maintaining information in WM should systematically
140 reflect the memory content⁴. We tested this using an impulse stimulus to ‘ping’ potentially
141 hidden neural states (Fig. 2e). As predicted, the impulse-specific response clearly differentiated
142 the content of WM (Fig. 2f), even though the driving input (‘ping’) was held constant on each
143 trial. The decodability of the cued item showed a significant cluster from 148 to 398 ms after
144 impulse stimulus onset (permutation test, $n = 30$, $p = 0.002$, corrected, cluster-forming threshold
145 $p < 0.05$). Average decodability from 100 to 500 ms was also significant ($p = 0.004$). Cued item
146 decoding was also higher than task-irrelevant (uncued) item decoding (cluster: 216 to 386 ms, p
147 $= 0.009$, corrected; average: $p = 0.028$). Indeed, the uncued item showed no evidence for
148 decoding (no corrected clusters; average: $p = 0.687$), suggesting that content can be rapidly
149 purged from WM when instructed, leaving effectively no trace in the neural state.

150 To test whether the impulse response reflects a literal ‘reactivation’ of item-specific activity
151 observed during encoding (e.g., Fig. 2b), we also examined whether a classifier trained on the
152 activity elicited by the memory stimuli during encoding could be used to decode the memory
153 item during the impulse epoch (and vice versa). However, we found no evidence for significant
154 cross-generalization between discriminative activity patterns during encoding and discriminative
155 activity driven by the impulse (corrected clusters, $p > 0.347$). We propose that the impulse
156 stimulus simply acts as a functional ping to recover hidden states, rather than a literal
157 ‘reactivation’ of a latent representation²¹.

158

159 Trial-wise variability in decoding the impulse response also predicted variability in WM
160 performance. Higher decoding trials of the cued item were accompanied by higher performance
161 than low decoding trials (permutation test, $n = 30$, $p = 0.043$; Fig 3a, left). There was also a
162 complementary cost for decoding the uncued item (i.e., a high decoding score for the uncued
163 item led to a decrease in accuracy on the cued item; $p = 0.002$; Fig 3b, left), suggesting that

164 participants might have failed to discard the uncued item (or simply did not use the cue properly)
 165 on some trials, contributing to error in performance. Finally, the difference between the accuracy
 166 effect of the cued and the uncued item was also significant (permutation test, $n = 30$, $p < 0.001$).



167

168 **Figure 3. Relationship between item-specific impulse decoding and WM accuracy. a.**

169 Difference in overall WM task performance between high and low cued item decoding trials
 170 (left). Proportion clockwise response for high and low decoding trials as a function of the angular
 171 difference between the memory item and the probe (right). Inset shows the difference in the
 172 slope parameter (a measure of memory precision) between high and low decoding trials. Data
 173 points outside of the 1.5 * interquartile range are shown separately in the boxplots (small
 174 crosses). Superimposed circles and error-bars are the mean and 95% C. I. of the mean. **b.** The
 175 same convention as in a. but for the decoding of the uncued item. Significant differences in

176 accuracy/precision between high and low decoding trials are highlighted by asterisks
177 (permutation test, $n = 30$, $p < 0.05$).

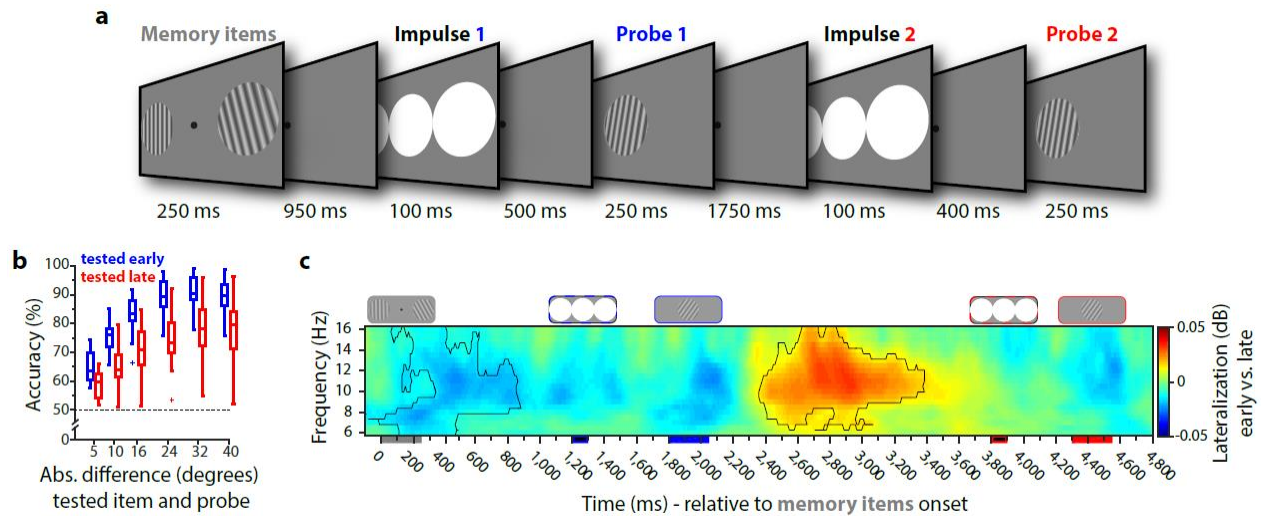
178 In principle, the relationship between trial-wise decoding and WM performance may rest on an
179 increase in guess-rate (i.e., due to forgetting or failure to encode), or a reduction in precision, or
180 both^{28,29}. To separate these possible contributions, we modelled the behavioural profile over
181 degrees of angular rotation between the memory item and the probe stimulus (see Online
182 Methods)^{30,31}. We found that the link to behaviour is most likely driven by a decrease in
183 precision (the slope parameter of the model) for weakly encoded hidden states of WM
184 (permutation test, $n = 30$, $p = 0.023$, one-tailed; Fig. 3a, right), while no evidence for an effect in
185 guess rate (the asymptote parameter) was found ($p = 0.867$, one-tailed). Modelling the observed
186 uncued item accuracy effect was inconclusive (Fig. 3b, right), with no evidence for either a
187 precision or guess rate effect ($p = 0.443$ and $p = 0.184$ respectively, one-tailed). Finally, we
188 found no evidence that trial-wise item decoding during the initial presentation of the memory
189 stimuli relates to memory performance (Supplementary Fig. 2a), further suggesting that the
190 relationship between accuracy and decoding triggered by the impulse is not due to a failure to
191 encode the memory item.

192

193 Experiment 2

194 Recently, it has been proposed that information in WM can be represented in qualitatively
195 different states^{32–34}, with attended items encoded in activity states measurable with standard
196 recordings of delay activity, whereas activity-silent states could underlie the representation of
197 currently unattended information in WM. In Experiment 2 ($n = 19$) we test whether unattended
198 but nevertheless remembered information in WM can still be decoded from the impulse
199 response. Again, two memory items were presented at the start of the trial, however both were
200 ultimately relevant as they would both be probed. Priority was manipulated by blocking the order
201 in which items would be probed (Fig. 4a), and instructing participants accordingly. Because there
202 was no other clue as to which item was being probed first or second, non-random responses
203 already indicate that participants used this blocked information (Fig. 4b). This was further
204 supported by lateralised changes in alpha power (Fig. 4c). During and shortly after the initial
205 presentation of the memory stimuli, there was a relative decrease in power at sensors
206 contralateral to the initially prioritised item, consistent with selective allocation of attention

207 (permutation test, $n = 19$, $p = 0.023$, corrected, cluster-forming threshold $p < 0.05$). Moreover,
 208 this pattern reversed after the response to the first item ($p = 0.009$, corrected), consistent with the
 209 assumption that participants then shift the originally de-prioritised item into the focus of
 210 attention in WM in preparation for the second probe ³⁵.



211

212 **Figure 4. Experiment 2 task structure, behavioural performance and attention-related**

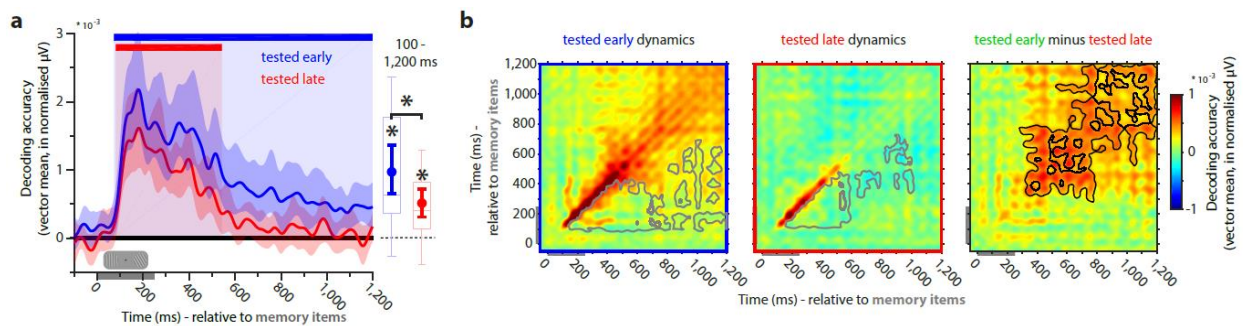
213 **alpha band activity.** **a.** Trial schematic. Two memory items were presented. Participants were
 214 instructed to maintain both items and were told at the start of each block which order the items
 215 would be tested. The first impulse was presented within the first memory delay (maintain both
 216 items, but attend the prioritised item), after which the prioritised item was probed. The second
 217 impulse was presented during the subsequent memory delay (maintain and attend only the now-
 218 prioritised item), after which the remaining item was probed. **b.** Boxplots show the accuracy of
 219 the early and late tested item as a function of the absolute angular difference (in degrees)
 220 between the memory item and the probe. Data points outside of the $1.5 \times$ interquartile range are
 221 shown separately in the boxplots (small crosses). **c.** Time-frequency representation of the
 222 difference between the contra- and ipsilateral posterior electrodes relative to the presentation side
 223 of the early tested memory items. Highlighted areas indicate significant difference (permutation
 224 test, $n = 19$, cluster-forming threshold $p < 0.05$, corrected significance level $p < 0.05$).

225

226 Decoding during stimulus presentation: We first analysed decoding during the initial processing
 227 of the memory stimuli. The results are plotted separately as a function of test-time (early or late

228 in the trial) as this could be meaningfully classified from the beginning of the trial (Fig. 5a). As
 229 expected, decoding the prioritised item (cluster: 74 to 1,200 ms, $p < 0.001$, corrected, cluster-
 230 forming threshold $p < 0.05$; average: $p < 0.001$), relative to the de-prioritised item (cluster: 82 to
 231 542 ms, corrected, $p < 0.001$, corrected; average: $p < 0.001$) was more robust (average: $p =$
 232 0.013). While decoding of the unattended item drops to chance relatively quickly after item
 233 presentation, the attended item shows significant decoding until the end of the epoch, replicating
 234 previous evidence showing that maintenance of only attended WM items is represented in the
 235 recorded brain activity patterns^{9–11}.

236



237

238 **Figure 5. Priority-dependent encoding and maintenance in WM.** **a.** Decodability of the item
 239 that is tested early (blue) and the item that is tested late (red) during memory item presentation.
 240 Blue and red bars indicate significant decoding clusters for the early- and late-tested item,
 241 respectively (permutation test, $n = 19$, cluster-defining threshold $p < 0.05$, corrected significance
 242 level $p < 0.05$). Error shading is 95% C.I. of the mean. Boxplots and superimposed circles with
 243 error bars (mean and 95 % C.I. of the mean) represent average decodability from 100 ms after
 244 stimulus onset until the end of the epoch. Significant average decoding and average difference
 245 between the decodability of the early and late item are marked by an asterisk (permutation test, n
 246 = 19, $p < 0.05$). **b.** Cross-temporal decoding matrices of the early (left) and late-tested (middle)
 247 item derived from training and testing on all time-point combinations, and the difference
 248 between the decoding of the early and late tested item (right). The grey outline indicates time-
 249 points of significantly lower decoding relative to both equivalent time-points along the diagonal,
 250 which is taken as evidence for dynamic coding (permutation test, $n = 19$, cluster-defining
 251 threshold $p < 0.05$, corrected significance level $p < 0.05$). The black outline (right) indicates

252 significantly higher decodability of the early compared to the late tested item (permutation test, n
253 = 19, cluster-defining threshold $p < 0.05$, corrected significance level $p < 0.05$).

254 The difference between attended and unattended item-maintenance in WM was even more
255 apparent when comparing their cross-temporal decoding matrices. Minimal cross-temporal
256 generalization during and shortly after memory item presentation suggested highly dynamic item
257 encoding: orientation discriminative patterns change over time. This was supported by
258 significant dynamic coding clusters during item encoding for both the early and late tested item,
259 where off-diagonal time-points show significantly lower decodability than both corresponding
260 on-diagonal time-points (permutation test, $n = 19$, cluster-defining threshold $p < 0.05$, corrected
261 significance level $p < 0.05$; see Online Methods; Fig. 5b, left and middle). However, the attended
262 item clearly showed a more time-invariant decoding pattern at the end of the epoch than the
263 unattended item, apparent by both significantly higher decodability on same time-point as well as
264 cross time-point decoding ($n = 19$, $p = 0.023$, corrected, cluster-forming threshold $p < 0.05$; Fig.
265 5b right). This further suggests that while the attended item also has a corresponding WM
266 maintenance signature in stable activity patterns, the unattended item does not.

267

268 Decoding of the impulse responses: Critically, we found that both the attended (clusters: 80 to
269 308 ms, $p = 0.004$, and 332 ms to 434 ms, $p = 0.031$, corrected; average: $p < 0.001$) and
270 unattended items (cluster: 172 to 306 ms, $p = 0.011$, corrected; average: $p = 0.045$) were
271 decodable in the first impulse response (Fig. 6a). This contrasts with the clear cueing differences
272 observed in Experiment 1, and suggests multiple items can be encoded in hidden states and
273 revealed by the impulse, even if only one item is in the focus of attention. It is worth noting,
274 however, that the decodability of the attended item was significantly higher than that of the
275 unattended item (average: $p = 0.031$), consistent with the behavioural evidence for relatively
276 better memory for the initially prioritised item.

291 the boxplot and error-bar of the difference in the slope parameter (a measure of memory
292 precision) between high and low decoding trials. **d.** The same convention as in a. but for the
293 decoding of the late-tested item during the late impulse. Significant differences in
294 accuracy/precision between high and low decoding trials are highlighted by asterisks
295 (permutation test, $n = 19$, $p < 0.05$, two-sided and one-sided for accuracy and precision tests,
296 respectively). Data points outside of 1.5 * interquartile range are shown separately in the
297 boxplots (small crosses).

298

299

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

304 We also found that the remaining relevant and initially unattended item could also be decoded in
305 the second impulse response (cluster: 196 to 326 ms, $p = 0.016$, corrected; average: $p = 0.012$),
306 while decoding the initially prioritised item failed to reach significance in this epoch (clusters: p
307 > 0.109 , corrected; average: $p = 0.112$; Fig 6b). The now-deprioritised item was presumably
308 cleared from the hidden state because it was no longer relevant, similar to forgetting observed
309 after the retro-cue from Experiment 1.

310 Again, we also tested for cross-generalization between the decodable patterns of the memory
311 items epoch (Fig. 5a) and the impulse-epochs (Fig. 6a, b). However, like in Experiment 1, we
312 found no evidence that the impulse literally ‘reactivates’ activity patterns associated with initial
313 encoding for either item (all corrected clusters: $p > 0.32$).

314 There was also a positive relationship between trial-wise decoding of the attended items at the
315 first and at the second impulse with WM performance (early: $p = 0.038$, Fig. 6c; late: $p = 0.04$;
316 Fig. 6d), replicating and extending the findings of Experiment 1. As in Experiment 1, we
317 modelled the behavioural profile to test if the positive relationship between decoding and task
318 performance is due to an increase in precision and/or a decrease in guess-rate. While the
319 modelling results were inconclusive for the early-tested item (precision: $p = 0.399$, one-tailed;

320 guess-rate: $p = 0.329$, one-tailed; Fig. 6c), there was evidence for an effect in precision of
321 working memory for the late item (precision: $p = 0.006$, one-tailed; guess-rate: $p = 0.942$, one-
322 tailed; Fig. 6d), replicating the precision effect of Experiment 1. Note that there was again no
323 relationship between accuracy and item decoding during the encoding phase (Supplementary Fig.
324 2b).

325 Experiment 3

326 We developed the impulse perturbation approach to reveal otherwise hidden neural states,
327 without necessarily transforming the mnemonic representation^{4,21}. This contrasts with other
328 studies using retro-cues^{10,11,32} or TMS³⁶ to ‘reactivate’ a latent item in working memory.
329 However, to test whether our impulse stimulus actually did result in a behaviourally relevant
330 transformation of the memory item (i.e., from a functionally latent to active state), we conducted
331 an additional behavioural experiment ($n = 20$). Adapting the design of Experiment 1, we now
332 varied the presentation of the stimulus-onset asynchrony (SOA) between impulse and probe
333 onset in Experiment 3 (SOA from 0 to 500ms; see Supplementary Fig. 4a). If the increase in
334 impulse-specific decodability observed in both EEG experiments reflects a functional
335 “reactivation” of an otherwise latent memory item, there should be a corresponding benefit to
336 behaviour.

337 A repeated measures ANOVA provided no evidence for an effect of SOA ($F(4, 76) = 1.184, p =$
338 0.325). Uncorrected paired comparisons between the no-impulse condition (SOA 0 ms) and all
339 other SOAs provided no evidence for an impulse-specific effect on accuracy for any SOA either
340 (permutation test, $n = 20$, all $p > 0.12$; Supplementary Fig. 4b). This suggests that our impulse
341 stimulus is effective for ‘pinging’ activity silent neural states, without resulting in any
342 behaviourally relevant transformation of the mnemonic representation.

Discussion

343

344 Recent theoretical models of WM predict a key role for activity-silent neural states in
345 maintaining item-specific information ^{4,17,18}. This raises a particular challenge for contemporary
346 neuroscience that is dominated by measurement and analysis of neural activation states. Here, we
347 address this challenge using a perturbation approach to reveal hidden neural states that code the
348 contents of WM. We show that the response to an impulse stimulus faithfully reflects item-
349 specific information in WM. We further demonstrate that the impulse response reflects both
350 attended and unattended items in WM, yet recently forgotten information leaves no detectable
351 trace in the hidden state. Behavioural modelling further suggests that the hidden-state coding
352 determines the quality of information in WM.

353 Previous evidence from non-human primates showed that a neutral visual stimulus presented
354 during the WM delay period can elicit distinct patterns of neural activity that depend on recent
355 visual input ³⁷. Although the previous work could not deconfound previous sensory stimulation
356 and WM proper, the observed effect helped motivate a dynamic coding model for WM ⁴.
357 According to this framework, distinct memoranda are associated with distinct changes in neural
358 response profile, which would be readable to downstream systems from the state-dependent
359 response to a retrieval probe ^{4,18}. Crucially, WM depends on the maintenance of the item-specific
360 neural response profile, rather than an explicit representation of an item in a persistent activity
361 state. We now provide direct evidence for a WM-dependent impulse response decoupled from
362 previous stimulation history, and further demonstrate that this WM state is highly flexible and
363 coupled to behavioural performance. The hidden state for a specific item can be rapidly cleared if
364 it is no longer relevant to the task, providing a striking neural correlate of directed forgetting in
365 WM.

366 Recent retro-cuing evidence suggests that prioritising one WM item relative to other task-
367 relevant items improves neural decoding of the cued item, whereas decoding of unattended items
368 drops to chance levels even though the unattended information is still ultimately task relevant
369 and retrievable at the end of the trial ¹⁰. Item-specific delay activity therefore seems to reflect the
370 focus of attention, rather than WM per se ³². The impulse response reported here clearly differs
371 from the typical profile observed for decoding delay activity patterns. In Experiment 2, both
372 attended and unattended items could be decoded from the impulse response of the hidden state as

373 long as they are both still ultimately required for task performance. This suggests that if the
374 information is successfully maintained in WM, there is a corresponding trace in the hidden state,
375 irrespective of attentional priority. These results highlight the flexibility of WM, independently
376 of switching attention between specific items in WM. Activity states appear to track the focus of
377 attention^{10,11,32}, whereas hidden states, as revealed by the impulse response, more closely track
378 the actual contents of WM.

379 Exactly how the proposed hidden state can be used for WM-guided behaviour remains an
380 important open question. Computationally, supervised learning could determine the mapping
381 between the memory-dependent probe response and the correct behavioural response³⁸, however
382 such a learning strategy seems implausible for real-world behaviour. Trial and error learning of
383 arbitrary patterns does not seem a realistic model for WM, at least for humans. Instead, the
384 inherent dynamics could establish a history-dependent match filter²⁰, which would be capable of
385 transforming probe input to a common decision signal (i.e., match/no-match, or in our case
386 clockwise/counter-clockwise). In Myers et al.²⁷, such a mechanism was shown to generate two
387 distinct decision-related signals in an orientation detection task: a signed (i.e., directional) and
388 unsigned difference signal, even though the signed difference was actually irrelevant to
389 behaviour in that task. A similar process could underpin WM encoding in hidden states. The
390 hidden state could establish a flexible, task dependent circuit for WM-dependent decision-
391 making³⁹. When the probe stimulus is presented, the hidden state transforms the input to
392 decision-relevant output: e.g., direction of angular rotation. However, because the impulse
393 stimulus used in these experiments does not contain decision-relevant features, the impulse
394 response reflects an input-output transformation of the arbitrary input.

395 It may be noted that although the response to an arbitrary input is sufficient to ‘read-out’ the
396 hidden state, it is unlikely to constitute an explicit ‘reactivation’ of the memory representation. In
397 contrast, retro-cueing can convert an unattended item to a prioritised state in preparation for the
398 recall²². Similarly, a recent transcranial magnetic stimulation study suggests that stimulation of
399 the visual cortex can also render an item active from its latent state³⁶. We find no evidence that
400 our impulse stimulus reactivates the same pattern associated with stimulus processing. Moreover,
401 a further behavioural experiment designed to test the possible behavioural consequences of our
402 impulse stimulus provides no evidence that it interacts with the mnemonic representation.
403 Rather, we argue that the impulse response simply ‘echoes’ the representational structure of the

404 hidden state, but does not drive an explicit transformation of latent memories to a prioritized
405 state.

406 It has long been assumed that WM maintenance depends on persistent neural activity ². Instead,
407 we propose that activity-silent neural states are sufficient to bridge memory delays. Activity-
408 dependent transformations in hidden states determine the temporary coding properties of
409 memory networks: i.e., dynamic coding ^{4,37}. WM decisions are made by the state-dependent
410 response to subsequent input. However, WM is also classically associated with active
411 manipulation of content in short-term memory ¹. We argue that such transformations are activity
412 dependent, but the results of the transformation can be maintained in short term memory via
413 latent network states. This alternative account does not ignore previous evidence for decodable
414 activity during mnemonic delays, but rather attributes such evidence to focused attention ³⁶,
415 periodic ¹⁸ or stochastic ¹⁷ updating, and/or response preparation ⁸. Interestingly, our current
416 results also show that cue-directed forgetting can rapidly wipe the mnemonic representation from
417 the hidden state. Rapid construction and dissolution of hidden states places important constraints
418 on the basic mechanisms of hidden-state coding.

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

429

Acknowledgements

430 We would like to thank the Biotechnology & Biological Sciences Research Council
431 (BB/M010732/1 to M.G.S.), and the National Institute for Health Research Oxford Biomedical
432 Research Centre Programme based at the Oxford University Hospitals Trust, Oxford University.
433 The views expressed are those of the authors and not necessarily those of the NHS, the NIHR, or
434 the Department of Health. We would also like to thank E. Spaak, A. Cravo, and N. Myers for
435 helpful comments, and all our volunteers for their participation.

436

437

Author Contributions

438 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.
439 analysed the data. M.J.W., M.G.S., E.G.A., and J.J. wrote the manuscript.

440

441

Competing Financial Interests

442 The authors declare no competing financial interests.

- 443 1. Baddeley, A. Working memory: looking back and looking forward. *Nat. Rev. Neurosci.* **4**,
444 829–839 (2003).
- 445 2. Curtis, C. E. & D’Esposito, M. Persistent activity in the prefrontal cortex during working
446 memory. *Trends Cogn. Sci.* **7**, 415–423 (2003).
- 447 3. Goldman-Rakic, P. Cellular basis of working memory. *Neuron* **14**, 477–485 (1995).
- 448 4. Stokes, M. G. ‘Activity-silent’ working memory in prefrontal cortex: a dynamic coding
449 framework. *Trends Cogn. Sci.* **19**, 394–405 (2015).
- 450 5. Watanabe, K. & Funahashi, S. Neural mechanisms of dual-task interference and cognitive
451 capacity limitation in the prefrontal cortex. *Nat. Neurosci.* **17**, 601–611 (2014).
- 452 6. Watanabe, K. & Funahashi, S. Prefrontal Delay-Period Activity Reflects the Decision
453 Process of a Saccade Direction during a Free-Choice ODR Task. *Cereb. Cortex* **17**, i88–i100
454 (2007).
- 455 7. Miller, E. K., Erickson, C. A. & Desimone, R. Neural Mechanisms of Visual Working
456 Memory in Prefrontal Cortex of the Macaque. *J. Neurosci.* **16**, 5154–5167 (1996).
- 457 8. Barak, O., Tsodyks, M. & Romo, R. Neuronal Population Coding of Parametric Working
458 Memory. *J. Neurosci.* **30**, 9424–9430 (2010).
- 459 9. LaRocque, J. J., Lewis-Peacock, J. A., Drysdale, A. T., Oberauer, K. & Postle, B. R.
460 Decoding Attended Information in Short-term Memory: An EEG Study. *J. Cogn. Neurosci.*
461 **25**, 127–142 (2012).
- 462 10. Lewis-Peacock, J. A., Drysdale, A. T., Oberauer, K. & Postle, B. R. Neural Evidence for a
463 Distinction between Short-term Memory and the Focus of Attention. *J. Cogn. Neurosci.* **24**,
464 61–79 (2011).

- 465 11. Sprague, T. C., Ester, E. F. & Serences, J. T. Restoring Latent Visual Working Memory
466 Representations in Human Cortex. *Neuron* **91**, 694–707 (2016).
- 467 12. Sreenivasan, K. K., Curtis, C. E. & D’Esposito, M. Revisiting the role of persistent neural
468 activity during working memory. *Trends Cogn. Sci.* **18**, 82–89 (2014).
- 469 13. Buonomano, D. V. & Maass, W. State-dependent computations: spatiotemporal processing
470 in cortical networks. *Nat. Rev. Neurosci.* **10**, 113–125 (2009).
- 471 14. Barak, O. & Tsodyks, M. Working models of working memory. *Curr. Opin. Neurobiol.* **25**,
472 20–24 (2014).
- 473 15. Murray, J. D. *et al.* Stable population coding for working memory coexists with
474 heterogeneous neural dynamics in prefrontal cortex. *Proc. Natl. Acad. Sci.* **114**, 394–399
475 (2017).
- 476 16. Fujisawa, S., Amarasingham, A., Harrison, M. T. & Buzsáki, G. Behavior-dependent short-
477 term assembly dynamics in the medial prefrontal cortex. *Nat. Neurosci.* **11**, 823–833 (2008).
- 478 17. Lundqvist, M. *et al.* Gamma and Beta Bursts Underlie Working Memory. *Neuron* **90**, 152–
479 164 (2016).
- 480 18. Mongillo, G., Barak, O. & Tsodyks, M. Synaptic Theory of Working Memory. *Science* **319**,
481 1543–1546 (2008).
- 482 19. Hempel, C. M., Hartman, K. H., Wang, X.-J., Turrigiano, G. G. & Nelson, S. B. Multiple
483 Forms of Short-Term Plasticity at Excitatory Synapses in Rat Medial Prefrontal Cortex. *J.*
484 *Neurophysiol.* **83**, 3031–3041 (2000).
- 485 20. Sugase-Miyamoto, Y., Liu, Z., Wiener, M. C., Optican, L. M. & Richmond, B. J. Short-Term
486 Memory Trace in Rapidly Adapting Synapses of Inferior Temporal Cortex. *PLOS Comput*
487 *Biol* **4**, e1000073 (2008).

- 488 21. Wolff, M. J., Ding, J., Myers, N. E. & Stokes, M. G. Revealing hidden states in visual
489 working memory using electroencephalography. *Front. Syst. Neurosci.* **9**, (2015).
- 490 22. Griffin, I. C. & Nobre, A. C. Orienting Attention to Locations in Internal Representations. *J.*
491 *Cogn. Neurosci.* **15**, 1176–1194 (2003).
- 492 23. Landman, R., Spekreijse, H. & Lamme, V. A. F. Large capacity storage of integrated objects
493 before change blindness. *Vision Res.* **43**, 149–164 (2003).
- 494 24. Harrison, S. A. & Tong, F. Decoding reveals the contents of visual working memory in early
495 visual areas. *Nature* **458**, 632–635 (2009).
- 496 25. Worden, M. S., Foxe, J. J., Wang, N. & Simpson, G. V. Anticipatory biasing of visuospatial
497 attention indexed by retinotopically specific alpha-band electroencephalography increases
498 over occipital cortex. *J. Neurosci. Off. J. Soc. Neurosci.* **20**, (2000).
- 499 26. Saproo, S. & Serences, J. T. Spatial Attention Improves the Quality of Population Codes in
500 Human Visual Cortex. *J. Neurophysiol.* **104**, 885–895 (2010).
- 501 27. Myers, N. E. *et al.* Testing sensory evidence against mnemonic templates. *eLife* **4**, e09000
502 (2015).
- 503 28. Zhang, W. & Luck, S. J. Discrete fixed-resolution representations in visual working memory.
504 *Nature* **453**, 233–235 (2008).
- 505 29. Bays, P. M. & Husain, M. Dynamic Shifts of Limited Working Memory Resources in
506 Human Vision. *Science* **321**, 851–854 (2008).
- 507 30. Murray, A. M., Nobre, A. C. & Stokes, M. G. Markers of preparatory attention predict visual
508 short-term memory performance. *Neuropsychologia* **49**, 1458–1465 (2011).
- 509 31. Prins, N. & Kingdom, F. A. A. Palamedes: Matlab routines for analyzing psychophysical
510 data. <http://www.palamedestoolbox.org> (2009).

- 511 32. Larocque, J. J., Lewis-Peacock, J. A. & Postle, B. R. Multiple neural states of representation
512 in short-term memory? It's a matter of attention. *Front. Hum. Neurosci.* **8**, 5 (2014).
- 513 33. Olivers, C. N. L., Peters, J., Houtkamp, R. & Roelfsema, P. R. Different states in visual
514 working memory: when it guides attention and when it does not. *Trends Cogn. Sci.* (2011).
515 doi:10.1016/j.tics.2011.05.004
- 516 34. Souza, A. S. & Oberauer, K. In search of the focus of attention in working memory: 13 years
517 of the retro-cue effect. *Atten. Percept. Psychophys.* 1–22 (2016). doi:10.3758/s13414-016-
518 1108-5
- 519 35. Ede, F. van, Niklaus, M. & Nobre, A. C. Temporal expectations guide dynamic prioritization
520 in visual working memory through attenuated alpha oscillations. *J. Neurosci.* 2272–16
521 (2016). doi:10.1523/JNEUROSCI.2272-16.2016
- 522 36. Rose, N. S. *et al.* Reactivation of latent working memories with transcranial magnetic
523 stimulation. *Science* **354**, 1136–1139 (2016).
- 524 37. Stokes, M. G. *et al.* Dynamic Coding for Cognitive Control in Prefrontal Cortex. *Neuron* **78**,
525 364–375 (2013).
- 526 38. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation
527 by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- 528 39. Martínez-García, M., Rolls, E. T., Deco, G. & Romo, R. Neural and computational
529 mechanisms of postponed decisions. *Proc. Natl. Acad. Sci.* **108**, 11626–11631 (2011).

Online Methods

530

531 Participants

532 Thirty healthy adults (13 female, mean age 24.9 years, range 18-38 years) were included in the
533 analyses of Experiment 1, 19 (10 female, mean age 24.7 years, range 18-39 years) in Experiment
534 2, and 20 in Experiment 3 (13 female, mean age 21, range 18-29 years). During data collection
535 and preprocessing, 4 additional participants of Experiment 1, 1 additional participant of
536 Experiment 2, and 6 additional participants of Experiment 3 were excluded from all analyses due
537 to either low average performance on the memory task (below 60% accuracy) or excessive eye-
538 movements (more than 30% of trials contaminated). No statistical methods were used to pre-
539 determine sample sizes but our sample sizes are similar to those reported in previous publications
540 ^{21,28}. All participants of Experiment 1 and 2 received monetary compensation of £10/h, and
541 participation in Experiment 3 contributed to course credits. All participants gave written
542 informed consent. Experiments 1 and 2 were approved by the Central University Research Ethics
543 Committee of the University of Oxford and Experiment 3 was approved by the Departmental
544 Ethical Committee of the University of Groningen.

545 Apparatus and Stimuli

546 The experimental stimuli were generated and controlled by Psychtoolbox ⁴⁰, a freely available
547 MATLAB extension. The stimuli were presented on a 23" screen running at 100 Hz and a
548 resolution of 1920 by 1080 in Experiment 1, on a 22" screen at a resolution of 1680 by 1050 in
549 Experiment 2, and on a 19 inch CRT screen running at 100 Hz and a resolution of 1280 by 1024
550 in Experiment 3. Viewing distance was set at 64 cm in Experiment 1, 67.5 cm in Experiment 2
551 and approximately 60 cm (not controlled) in Experiment 3, to ensure that the visual angles of
552 stimuli were the same across experiments even though the screen parameters were different. A
553 standard keyboard was used for response input by the participants.

554 All reported stimuli were the same in all experiments, unless explicitly mentioned otherwise. A
555 grey background (RGB = 128, 128, 128; 20.5 cd/m²; 28.6 cd/m² in Experiment 3) was
556 maintained throughout the experiments. A black fixation dot with a white outline (0.242°) was
557 presented in the centre of the screen throughout all trials. Memory items and memory were sine-
558 wave gratings presented at 20% contrast, with a diameter of 6.69° and spatial frequency of 0.65
559 cycles per degree. The phase was randomized within and across trials. The memory items were

560 presented at 6.69° eccentricity and for each trial the orientations were randomly selected without
561 replacement from a uniform distribution of orientations. The impulse stimulus was 3 adjacent
562 ‘bullseyes’ in Experiment 1. Each ‘bullseye’ was of the same size and spatial frequency as the
563 memory items. To reduce strain on the eyes, and to minimise forward masking in Experiment 3,
564 the impulse stimulus in Experiments 2 and 3 consisted of 3 adjacent white circles. In Experiment
565 1 and 2 the probes had the same contrast and spatial frequency as the memory items, and was
566 presented in the centre of the screen. In Experiment 3 the probe screen included a high contrast
567 black and white square-wave grating in the centre and two white lateralized circles on the outside
568 (the same location and size as the preceding lateral impulse circles). The angle differences
569 between a memory item and the corresponding memory probe were uniformly distributed across
570 7 angle differences in Experiment 1 ($\pm 3^\circ$, $\pm 7^\circ$, $\pm 12^\circ$, $\pm 18^\circ$, $\pm 25^\circ$, $\pm 33^\circ$, $\pm 42^\circ$), 6 angle differences
571 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
572 Experiment 3.

573 Procedure

574 Experiment 1: Participants completed a retro-cue visual working memory task. Each trial began
575 with the onset of a fixation dot at the centre of the screen. After 1000 ms, the memory item array
576 was shown for 250 ms, consisting of two randomly oriented low-contrast gratings left and right
577 of fixation. After a delay of 800 ms an arrow was shown for 200 ms in the centre of the screen,
578 pointing either to the left or to the right, and thus cueing which of the two previously presented
579 items would be tested. The number of left and right cued trials was equal and the order was
580 randomized for each participant. The impulse stimulus was presented for 100 ms, 900 ms after
581 the offset of the retro-cue. After another delay of 400 ms, the memory probe was shown for 250
582 ms. Participants were instructed to indicate if the orientation of the probe relative to the
583 orientation of the memory item was rotated clockwise by pressing the “m” key with the right
584 index finger, or counter-clockwise by pressing the “c” key with the left index finger. A high or
585 low frequency feedback tone was played after response, indicating if the answer was correct or
586 incorrect, respectively. The next trial started within 400 to 700 ms (determined randomly).
587 Participants completed 1344 trials in total, which lasted approximately 3 hours (including
588 breaks). Trial conditions were randomized across the whole session. See Figure 1a for a trial
589 schematic.

590 Experiment 2: Participants completed a visual working memory task where two items were
591 serially tested. The experiment began by instructing the participant which of the two memory
592 items would be tested early, and which one would be tested late. This rule never changed within
593 a session. Each trial began with the onset of a fixation dot at the centre of the screen. After 1000
594 ms, the memory item array was shown for 250 ms, consisting of two randomly oriented low-
595 contrast gratings left and right from fixation. After a delay of 950 ms, the first impulse was
596 presented for 100 ms. After a delay of 500 ms, the first memory probe was presented for 250 ms,
597 probing the first item. The response input was the same as in Experiment 1. After a fixed delay
598 of 1750 ms after the offset of the first probe, the second impulse was shown for 100 ms.
599 Following a delay of 400 ms, the second memory probe was presented for 250 ms, probing the
600 late-tested item. After the second response, two feedback tones were played, one for each
601 response, separately indicating whether the first and second answers were correct. Participants
602 completed two sessions of the task on two separate days, separated by approximately 1-2 weeks.
603 The testing order of the memory items was fixed within each sessions, and switched between
604 sessions (i.e. left item tested first in one session, right item tested first in the other session). The
605 order of the testing rule between sessions, (whether the left item would be tested first in the first
606 or in the second session) was counterbalanced across participants (odd numbered left first, even
607 numbered right first). Each session consisted of 864 trials, and lasted approximately 3 hours
608 including breaks. See Figure 3a for a trial schematic.

609 Experiment 3: The task was almost the same as Experiment 1, including the same timings of the
610 memory items, cue, probe and overall trial duration. The one key difference was the timing of
611 the impulse stimulus. While the delay between cue offset and probe onset was held constant at
612 1,400 ms across all trials (the same as in Experiment 1), the SOA between impulse and probe
613 onset was 0, 50, 100, 250 or 500 ms (determined pseudo randomly across the session). No
614 impulse was shown in the 0 ms SOA condition. The impulse remained on the screen until the
615 probe stimulus was presented. This was to ensure the least possible interference of the impulse
616 on probe processing (i.e., rapid onset and offset of the white circles immediately before probe
617 presentation could deteriorate probe visibility), as well as keeping the different SOA conditions
618 as similar as possible (longer SOA would include an additional offset). Participants completed
619 280 trials (approximately 30 minutes). See Supplementary Fig. 4a for a trial schematic.

620 Data collection and analyses were not performed blind to the conditions of the experiments.

621 Due to the within-subject design in all three experiments, randomization of conditions between
622 subjects was not applicable.

623 EEG Acquisition

624 The EEG signal was acquired from 61 Ag/AgCl sintered electrodes (EasyCap, Herrsching,
625 Germany) laid out according to the extended international 10-20 system. Data was recorded at
626 1000 Hz using NeuroScan SynAmps RT amplifier and Scan 4.5 software in Experiment 1 and
627 Curry 7 software in Experiment 2 (Compumedics NeuroScan, Charlotte, NC). The anterior
628 midline frontal electrodes (AFz) served as the ground. Bipolar electrooculography (EOG) was
629 recorded from electrodes placed above and below the right eye, and from electrodes placed to the
630 left of the left eye and to the right of the right eye. The impedances of all electrodes were kept
631 below 5 k Ω . Online, the EEG was referenced to the right mastoid and filtered using a 200 Hz
632 low-pass filter.

633 EEG pre-processing

634 Offline, the data was re-referenced to the average of both mastoids, down-sampled to 500 Hz and
635 band-pass filtered (0.1 Hz high-pass and 40 Hz low-pass) using EEGLAB⁴¹. The data was then
636 epoched to the onset of the memory items and the impulse. In Experiment 1, the memory item
637 epoch was from -200 ms to 1050 ms, relative to onset, and in Experiment 2 from -200 ms to
638 1200 ms. The impulse epochs were from -200 ms to 500 ms relative to onset in both
639 experiments. Additionally, for the purpose of artefact rejection, which included the rejection of
640 trials containing saccadic eye-movement prior to the time of interest (see below), the cue
641 segment in Experiment 1 was also epoched (-200 ms to 1100 ms).

642 Subsequent artefact detection and trial rejection focused exclusively on the 17 posterior channels
643 that were included in the analyses (P7, P5, P3, P1, Pz, P4, P6, P8, PO7, PO3, POz, PO4, PO8,
644 O1, Oz, O2) and the EOGs. Each trial of each epoch was individually visually inspected for
645 blinks, saccades and non-stereotyped artefacts. Trials from individual epochs were rejected from
646 analyses involving that epoch if it contained any of the above-mentioned artefacts. Furthermore,
647 impulse-epoch trials were also excluded from corresponding analyses if the EOG signal
648 suggested that saccades occurred during any of the previous epochs of that trial. In Experiment 1
649 this exclusion procedure was applied to the cue-epoch as well. In Experiment 2, late impulse
650 trials were also excluded if no response was registered for the preceding probe. For the decoding

651 analyses, each epoch was baselined using the average signal from -200 ms to 0 ms before
652 stimulus onset. The multivariate data were also demeaned at each time-point by subtracting the
653 average voltage for all posterior channels included in the analyses.

654 Time-frequency decomposition and lateralization analysis

655 In order to explore alpha power (8-12 Hz) lateralization^{25,42}, the spectral power from 6 to 16 Hz
656 (in steps of 0.5 Hz) of the EEG signal was computed using Hanning tapers with time-windows of
657 5 cycles per frequency (in steps of 10 ms) using the MATLAB toolbox FieldTrip⁴³. We included
658 the whole experimental trial, ranging from 1000 ms before memory item onset until 1500 ms
659 after (second) probe onset (-1000 to 4150 ms relative to memory items in Experiment 1, and -
660 1000 to 5800 ms relative to memory items in Experiment 2). The power was log transformed,
661 and lateralization was computed by subtracting the average power of the ipsilateral posterior
662 electrodes from the average power of the contralateral posterior electrodes in relation to the cued
663 memory item in Experiment 1 and to the early-tested item in Experiment 2 (P7, P5, P3, P1, PO7,
664 PO3, O1 versus P8, P5, P6, P4, P2, PO8, PO4, O2).

665 Significant clusters of lateralization were determined using a cluster-corrected non-parametric
666 sign-permutation test⁴⁴. In both experiments, the whole trial was included in this analysis (-100
667 to 3150 ms relative to memory items onset in Experiment 1, and -100 to 4800 ms in Experiment
668 2).

669 Orientation decoding

670 To test whether the activity pattern of the posterior EEG channels of interest contained
671 orientation-specific activity, we used the Mahalanobis distance⁴⁵ to compute the trial-wise
672 distances between the full range of possible orientations, and quantify to what extent the
673 computed distances adhere to the parametric circular space of the orientations¹¹. This approach
674 is an extension of the pairwise distance approach we used before²¹ and is conceptually similar to
675 the population tuning curve model²⁶.

676 The left and right presented items were decoded separately and independently within each
677 participant and experimental session. All 17 posterior channels (see above) were used for all
678 decoding analyses. The procedure followed a leave-one-trial-out cross-validation approach to
679 compute the trial-wise decodability of the orientation of interest. The activity pattern of a single

680 test-trial at a particular time-point was compared to the pattern of all other trials at the same
681 time-point. These were averaged into 12 orientation bins relative to the orientation of the test-
682 trial, each containing trials with orientations within a range of 30° and centred around -75° , -60° ,
683 -45° , -30° , -15° , 0° , 15° , 30° , 45° , 60° , 75° , and 90° . The Mahalanobis distances between the
684 test-trial and each orientation bin was computed using the covariance estimated from all trials
685 excluding the test-trial using a shrinkage estimator⁴⁶. To simplify visualization and
686 interpretation, the 12 resulting distances were mean centred and the sign was reversed, resulting
687 in a visual representation of a tuning curve. Higher values correspond to greater relative
688 similarity between the test-trial and the averaged train-trials within a particular orientation bin,
689 and lower values correspond to greater dissimilarity.

690 Next, the vector means of the tuning curves were computed¹¹. First, the cosine of the centre of
691 each orientation bin (θ) was rescaled to the range -180 to 180 . It was then multiplied with the
692 corresponding sign-reversed distances ($d(\theta)$) before the mean of the resulting 12 values was
693 taken, which made up the decoding accuracy (da).

694 Equation 1: $da = \text{mean}(d(\theta) \cos(2\theta))$

695 A high value reflects evidence for orientation tuning: the difference between the test-trial and
696 train-trials with a similar orientation is smaller than between the test-trial and train-trials with
697 different orientations. This procedure was repeated for all trials and all time-points. See
698 Supplementary Information for the custom Matlab function used to decode orientations using
699 Mahalanobis distance.

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

702 Cluster-corrected sign-permutation significance tests were carried out within the memory items
703 epoch (0 to 1050 ms in Experiment 1, 0 to 1200 ms in Experiment 2) and impulse epochs
704 separately (0 to 500 ms in both experiments), in order to explore the significant decoding time-
705 course. Additionally, to assess the overall decodability within an epoch, the decoding values
706 were averaged over time (from 100 ms after stimulus onset until the end of the epoch) and then
707 submitted to a two-sided permutation test.

708 Relationship between behaviour and decoding

709 The trial-wise average decoding scores after memory items presentation (100 to 1050 ms in
 710 Experiment 1, and 100 to 1200 ms in Experiment 2) and impulse presentation (100 ms to 500
 711 ms) was median split. Non-response trials (to the early probe in Experiment 2) were excluded
 712 from this analysis. The average behavioural accuracies of high and low decoding trials were
 713 statistically compared using a two-sided permutation test.

714 Behavioural modelling

715 To further explore the relationship between WM task performance and trial-wise decoding, we
 716 modelled the behavioural performance as a function the difference in degrees between the
 717 orientation of the memory item and the probe using the following model that was fit to each
 718 participant separately³⁰.

$$719 \text{ Equation 2: } y = \lambda + \frac{(1-2\lambda)}{2} \times \operatorname{erfc}\left(\frac{-\beta}{\sqrt{2}}(x - \alpha)\right)$$

720 where *erfc* is the complementary Gaussian error function, λ is the asymptote, β is the slope and α
 721 is the threshold/bias parameter. The modelling fitting was performed using the Palamedes
 722 Matlab toolbox³¹. The asymptote represents the guess rate, where a higher value reflects a
 723 higher probability that no information about the probed item is maintained in WM, resulting in a
 724 higher probability for mistakes even when the angular difference between the probe and the
 725 memory item is large. The slope is interpreted as the memory precision, where a high precision
 726 reflects a relatively high proportion of correct responses at small degree rotations between the
 727 probe and memory item. The asymptote and slope parameters were both unconstrained across the
 728 high and low decoding conditions. A single bias parameter was used, which was included
 729 (instead of fixing it at 0) because cumulative-likelihood tests⁴⁷ showed better model fits for all
 730 cases (Experiment 1: $n = 30$, $\chi^2(30) = 135.978$, $p < 0.001$; Experiment 2, $n = 19$, early accuracy:
 731 $\chi^2(19) = 215.351$, $p < 0.001$; late accuracy: $\chi^2(19) = 33.69$, $p = 0.02$).

732 The unconstrained model parameters (slope and asymptote) were subsequently compared
 733 between high and low decoding trials. Since the behavioural modelling was carried out as a
 734 direct follow up to the average accuracy effects observed in both experiments (two-sided tests),
 735 we had clear expectations about the directionality of the effects. For the positive relationship
 736 between decoding and accuracy observed for the cued item in Experiment 1 and both tests in
 737 Experiment 2, we expected that decoding should have a negative relationship with guess-rate

738 (i.e., lower guess-rate for high decoding) and/or a positive relationship with precision (higher
739 precision for higher decoding), and vice versa for the negative accuracy effect of the uncued item
740 in Experiment 1. Therefore, all tests of model parameter comparisons between high and low
741 decoding trials were one-sided.

742 Cross-temporal decoding

743 We also explored the cross-temporal dynamics of stimulus processing and maintenance as a
744 function of item priority in Experiment 2, and cross-generalization between impulse and memory
745 presentation epochs in both experiments. The decoding approach was the same as described
746 above, except classifiers trained at each time point were tested at every other time point,
747 resulting in 2-dimensional cross-temporal decoding matrices⁴⁸.

748 If the decoding patterns are stationary, it should not matter whether train/test is performed using
749 the same time points. In contrast, decoding often appears dynamic: training and testing on the
750 same time-points results in higher decoding scores than training and testing on different time-
751 points (i.e., minimal cross-temporal generalization). We tested for this hallmark feature of
752 dynamic coding using a non-parametric test used previously²⁷. The decodability at each cross-
753 temporal time-point $t_{x,y}$ was compared to the pair of decodabilities at the corresponding within
754 time-points ($t_{x,x}$ and $t_{y,y}$) with two separate permutation tests. A significant difference in both was
755 taken as evidence for dynamic coding. Time-points of significant dynamic coding were corrected
756 for multiple comparisons using a two-dimensional cluster-based permutation test.

757 Significance testing

758 To determine statistical significance, we used the non-parametric sign-permutation test⁴⁴(with
759 one exception, see ANOVA below), which does not make assumptions about the underlying
760 distribution. Since the null hypotheses of all tests corresponded to no effect (i.e. no difference in
761 power lateralization, no difference in decodability, etc.), the sign of the data of each participant
762 was randomly flipped with a probability of 50% 50,000 times. The resulting distribution was
763 used to derive at the p -value of the null-hypothesis that the mean effect was equal to 0. All tests
764 were two-sided, unless otherwise stated.

765 For time-series and frequency data, the above procedure was repeated for each time-point and
766 frequency (when applicable). To correct for multiple comparisons over time and/or frequencies,

767 a cluster-based permutation test was subsequently used using 50,000 permutations (5,000 for
768 cross-temporal decoding, due to computer memory limitations), with a cluster-forming threshold
769 and cluster significance threshold of $p < 0.05$. Tests concerning the average of specific time-
770 windows (which includes decoding-behaviour relationships) were performed to test unique and
771 independent hypotheses, therefore no correction applied. The sample size for all tests in
772 Experiment 1 was $n = 30$, $n = 19$ in Experiment 2, and $n = 20$ in Experiment 3.

773 The 95 % confidence intervals of the error-bars were determined by bootstrapping from the
774 corresponding data 50,000 times.

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

779 A repeated measures ANOVA was used to analyse the behavioural data of Experiment 3. The
780 normality and equal variances assumptions were tested with the Shapiro-Wilk test of normality
781 and Mauchly's test of sphericity, respectively. Neither test provided evidence for assumption
782 violations of the data.

783 Data availability

784 The data that support the finding of this study are publically available at
785 <http://datasharedrive.blogspot.co.uk/2017/03/dynamic-hidden-states-underlying.html>. All
786 necessary task/condition information has been provided within a self-contained format, as
787 specified in the OECD Principles and Guidelines for Access to Research Data from Public
788 Funding⁴⁹.

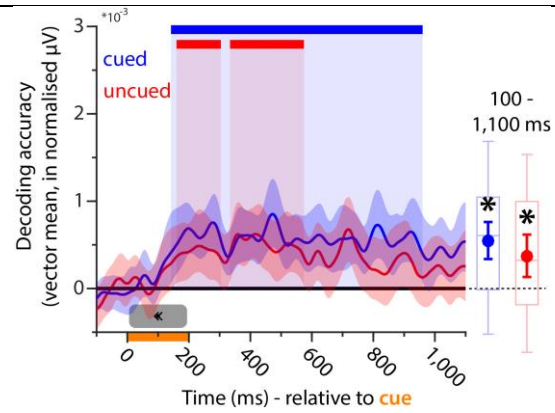
789 Code availability

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

793 **Methods References**

794 40. Brainard, D. H. The Psychophysics Toolbox. *Spat. Vis.* **10**, 433–436 (1997).

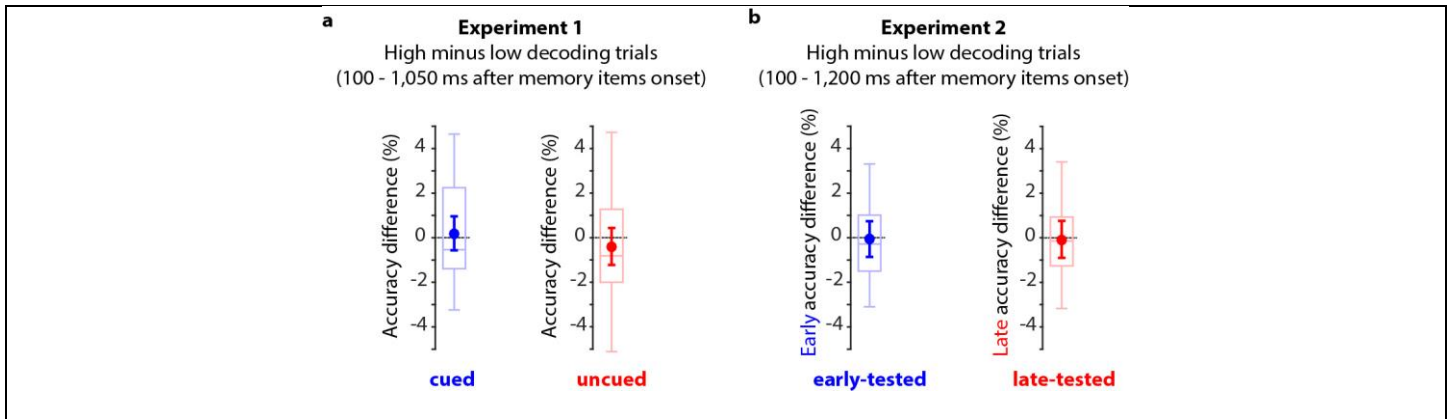
- 795 41. Delorme, A. & Makeig, S. EEGLAB: an open source toolbox for analysis of single-trial EEG
796 dynamics including independent component analysis. *J. Neurosci. Methods* **134**, 9–21
797 (2004).
- 798 42. Schneider, D., Mertes, C. & Wascher, E. The time course of visuo-spatial working memory
799 updating revealed by a retro-cuing paradigm. *Sci. Rep.* **6**, 21442 (2016).
- 800 43. Oostenveld, R. *et al.* FieldTrip: Open Source Software for Advanced Analysis of MEG,
801 EEG, and Invasive Electrophysiological Data, FieldTrip: Open Source Software for
802 Advanced Analysis of MEG, EEG, and Invasive Electrophysiological Data. *Comput. Intell.*
803 *Neurosci. Comput. Intell. Neurosci.* **2011**, **2011**, e156869 (2010).
- 804 44. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *J.*
805 *Neurosci. Methods* **164**, 177–190 (2007).
- 806 45. De Maesschalck, R., Jouan-Rimbaud, D. & Massart, D. L. The Mahalanobis distance.
807 *Chemom. Intell. Lab. Syst.* **50**, 1–18 (2000).
- 808 46. Ledoit, O. & Wolf, M. Honey, I shrunk the sample covariance matrix. *J. Portf. Manag.* **30**,
809 110–119 (2004).
- 810 47. Claessens, P. M. E. & Wagemans, J. A Bayesian framework for cue integration in
811 multistable grouping: Proximity, collinearity, and orientation priors in zigzag lattices. *J. Vis.*
812 **8**, 33–33 (2008).
- 813 48. King, J.-R. & Dehaene, S. Characterizing the dynamics of mental representations: the
814 temporal generalization method. *Trends Cogn. Sci.* **18**, 203–210 (2014).
- 815 49. Pilat, D. & Fukasaku, Y. OECD Principles and Guidelines for Access to Research Data from
816 Public Funding. *Data Sci. J.* **6**, OD4-OD11 (2007).
- 817



Supplementary Figure 1

Cue-specific item decoding time-course (Experiment 1).

The cue-specific neural response showed robust decoding for the cued ($n = 30$, cluster: 142 to 960 ms, $p < 0.001$, corrected; average: $p < 0.001$) and uncued item ($n = 30$, clusters: 158 to 304 ms, $p = 0.035$, corrected, 328 to 574 ms, $p = 0.006$, corrected, average: $p = 0.006$). Error shading is the 95 % C.I. of the mean. The boxplots and superimposed circles with error-bars (mean and 95 % C.I. of the mean) represent average decoding from 100 to 1,100 ms after cue onset. Significant average decoding is marked by an asterisk (permutation test, $n = 30$, $p < 0.05$).

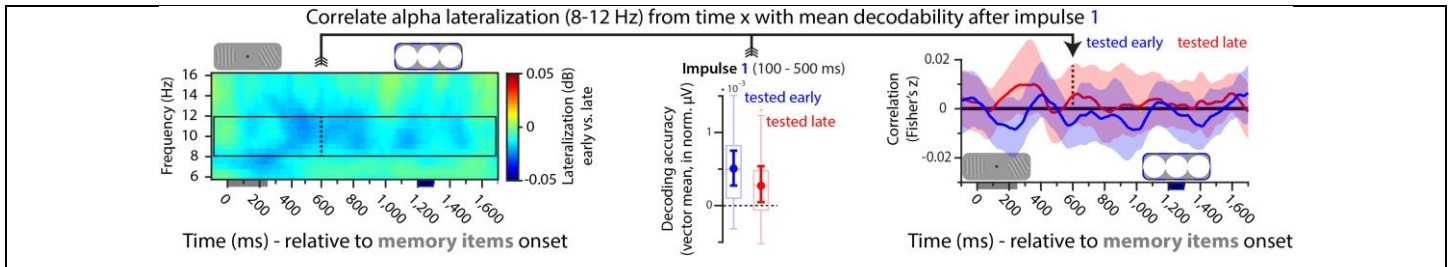


Supplementary Figure 2

Testing the relationship between item decoding during the memory item epoch and WM accuracy.

a. Task accuracy difference between high and low decoding trials of the cued (blue; $n = 30$, $p = 0.673$) and uncued (red; $n = 30$, $p = 0.344$) item during the memory items epoch (average decoding from 100 to 1,050 ms relative to memory items onset) in Experiment 1.

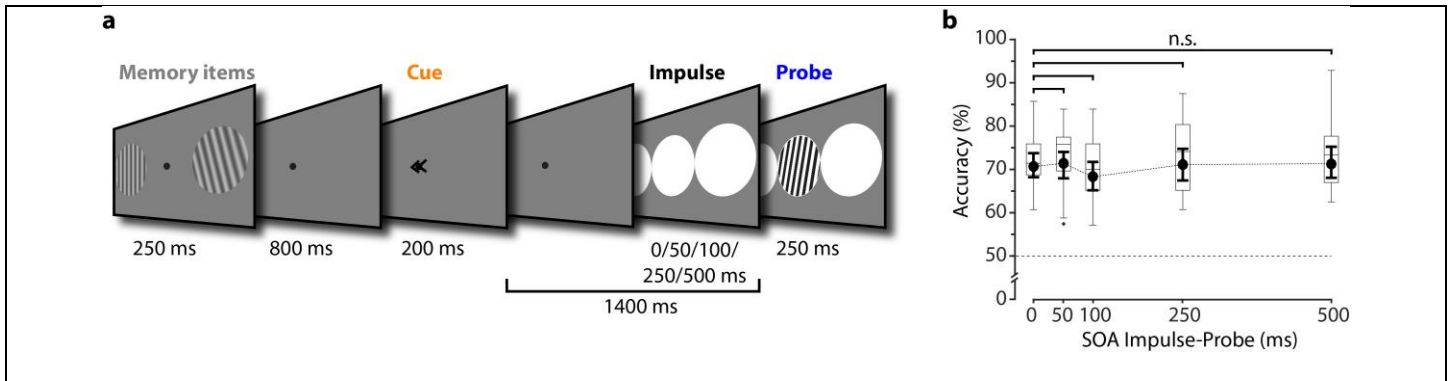
b. Early accuracy difference between high and low decoding trials of the early-tested item (blue; $n = 19$, $p = 0.865$), and late accuracy difference between high and low decoding trials of the late-tested (red; $n = 19$, $p = 0.978$) during the memory items epoch (average decoding from 100 to 1,200 ms relative to memory items onset) in Experiment 2. Circles and error-bars superimposed on the boxplots represent mean and 95% C.I. of the mean.



Supplementary Figure 3

Testing the relationship between alpha-lateralization and item decoding after the first impulse in Experiment 2.

Both attended and unattended memory items were decodable after the first impulse in Experiment 2; however, it remains possible that participants sometimes attended to the less-relevant item, contributing to decoding on some trials. To consider this possibility, we test whether the impulse-specific WM item decoding after impulse 1 presentation covaries with trial-wise fluctuations in spatial attention. Spatial attention was indexed by alpha-power lateralization relative to the location of the early-tested item of each time-point (left, also see Figure 4c and corresponding results), and trialwise item decodability was estimated 100-500ms after impulse 1 onset (middle panel). The correlation time-course (right), where each time-point represents the mean correlation of the averaged item decoding (100 – 500 ms after impulse 1) with the alpha-lateralization of that time-point, shows no evidence for a relationship between item decoding and alpha-lateralization for any time-point (permutation test, $n = 19$, early-tested item: all $p > 0.058$; late-tested item: all $p > 0.148$, uncorrected). Therefore, we find no evidence that the impulse-response varies with the focus of attention, even on a trial-wise basis. Error shadings are 95% C.I. of the mean. Circles and error bars superimposed on the boxplots represent mean and 95% C.I. of the mean. Data points outside the 1.5 * interquartile range are marked as crosses in the boxplots.



Supplementary Figure 4

Task schematic and results of behavioural experiment.

a. Two memory items were presented, and participants were instructed to memorize both. A retro-cue indicated which item would be tested at the end of the current trial. The impulse stimulus was presented at varying delays (or not at all) and stayed on screen until the probe was presented. Participants indicated whether the probe was rotated clockwise or anti-clockwise relative to the orientation of the cued item. **b.** Behavioural performance as a function of impulse-probe SOA. None of the uncorrected paired comparisons between the no-impulse condition (SOA 0 ms) and the other SOA conditions reached significance (permutation test, $n = 20$). Circles and error bars superimposed on the boxplots represent mean and 95% C.I. of the mean. Data points outside the 1.5 * interquartile range are marked as crosses in the boxplots.

819

820