

1 **Drifting codes within a stable coding scheme for** 2 **working memory**

3
4 Authors: Wolff, M. J.^{1,2}, Jochim, J.³, Akyürek, E. G.², Buschman, T. J.⁴, & Stokes, M. G.^{1,3}

- 5 1. Department Experimental Psychology, University of Oxford, Oxford, OX2 6GG,
6 United Kingdom
7 2. Department Experimental Psychology, University of Groningen, Groningen, 9712 TS,
8 The Netherlands
9 3. Oxford Centre for Human Brain Activity, University of Oxford, Oxford, OX3 7JX,
10 United Kingdom
11 4. Princeton Neuroscience Institute and Department of Psychology, Princeton
12 University, Princeton, NJ 08540, USA

13
14 Correspondence: michael.wolff@psy.ox.ac.uk, mark.stokes@psy.ox.ac.uk

15

Abstract

16 Working memory (WM) is important to maintain information over short time periods to
17 provide some stability in a constantly changing environment. However, brain activity is
18 inherently dynamic, raising a challenge for maintaining stable mental states. To investigate
19 the relationship between WM stability and neural dynamics, we used electroencephalography
20 to measure the neural response to impulse stimuli during a WM delay. Multivariate pattern
21 analysis revealed representations were both stable and dynamic: there was a clear difference in
22 neural states between time-specific impulse responses, reflecting dynamic changes, yet the
23 coding scheme for memorized orientations was stable. This suggests that a stable
24 subcomponent in WM enables stable maintenance within a dynamic system. A stable coding
25 scheme simplifies readout for WM-guided behaviour, whereas the low-dimensional dynamic
26 component could provide additional temporal information. Despite having a stable subspace,
27 WM is clearly not perfect – memory performance still degrades over time. Indeed, we find that
28 even within the stable coding scheme, memories drift during maintenance. When averaged
29 across trials, such drift contributes to the width of the error distribution.

30

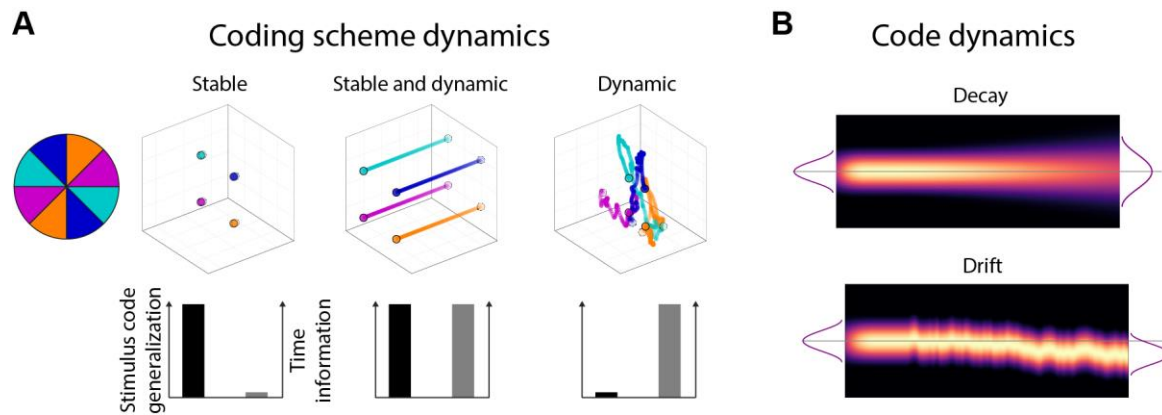
Introduction

31 Neural activity is highly dynamic, yet often we need to hold information in mind in a stable
32 state to guide ongoing behaviour. Working memory is a core cognitive function that provides
33 a stable platform for guiding behaviour according to time extended goals; however, it remains
34 unclear how such stable cognitive states emerge from a dynamic neural system.

35 At one extreme, WM could effectively pause the inherent dynamics by falling into a stable
36 attractor (e.g., 1,2). This solution has been well-studied, and provides a simple readout of
37 memory content irrespective of time (i.e., memory delay). However, more dynamic models
38 have also been suggested. For example, in a recent hybrid model, stable attractor dynamic
39 coexist with a low-dimensional, time varying component (3,4); see Fig. 1A for model
40 schematics). This permits some dynamic activity, whilst also maintaining a fixed coding
41 relationship of WM content over time (5). As in the original stable attractor model, the coding
42 scheme is stable over time, permitting easy and unambiguous WM read out by downstream
43 systems, regardless of maintenance duration (6). Finally, it is also possible to maintain stable
44 information in a richer dynamical system (e.g., 7). Although the relationship between activity
45 pattern and memory content changes over time, the representational geometry could remain
46 relatively constant (5). Such dynamics emerge naturally in a recurrent network, and provide
47 rich information about the previous input, and elapsed time (8), but necessarily entail a more
48 complex readout strategy (time-specific decoders or a high-dimensional classifier that finds a
49 high-dimensional hyperplane that separates memory condition for all time points - (9)).

50 Although all models seek to account for stable WM representation, it is also important to note
51 that maintenance in WM is far from perfect. In particular, WM performance decreases over
52 time. (10), which could be ascribed to two different mechanisms (Fig. 1B). On the one hand,
53 the neural representation could simply degrade over time, either due to an overall decrease in
54 WM specific neural activity, or through a general broadening of the neural representation (11).
55 In this framework, the distribution of recall error reflects sampling from a broad underlying
56 distribution. On the other hand, the neural representation of WM content might gradually drift
57 along the feature dimension as a result of the accumulating effect of random shifts due to noise
58 (12). Even if the underlying neural representation remains sharp, variance in the mean over
59 trials results in a relative broad distribution of errors over trials.

60



61 **Figure 1.** Model predictions. (A) The relationship between the neural coding scheme of
62 orientations (colours) in WM over time. Left: A stable coding scheme within a stable neural
63 population. Middle: A stable coding scheme within a dynamic neural population. Right: A
64 dynamically changing coding scheme. (B) The fidelity of the population code in WM over
65 time. Top: The code decays and becomes less specific over time, leading to random errors
66 during read-out. Bottom: The code drifts along the feature dimension, leading to a still sharp,
67 but shifted code during read-out.

68 Computational modelling based on behavioural recall errors from WM tasks with varying set-
69 sizes and maintenance periods predict a drift for colours and orientations maintained in WM
70 (13,14). At the neural level, evidence for drift has been found in the neural population code in
71 monkey prefrontal cortex during a spatial WM task (15), where trial-wise shifts in the neural
72 tuning profile predicted if recall error was clockwise or counter-clockwise relative to the
73 correct location. Recently, a human fMRI study has found that delay activity reflected the probe
74 stimulus more when participants erroneously concluded that it matched the memory item (16),
75 which is consistent with the drift account.

76 Tracking these neural dynamics of non-spatial neural representations, which are not related to
77 spatial attention or motor planning, is not trivial in humans. Previously we found that the
78 presentation of a simple impulse stimulus (task-relevant visual input) presented during the
79 maintenance period of visual information in WM results in a neural response that reflects non-
80 spatial WM content (17,18). Here we extend this approach to track WM dynamics. In the
81 current study we developed a paradigm to test the stability (and/or dynamics) of WM neural
82 states and the consequence for readout by “pinging” the neural representation of orientations
83 at specific time-points during maintenance.

84 We found that the coding scheme remained stable during the maintenance period, even-though
85 maintenance time was coded in an additional low-dimensional axis. We furthermore found that
86 the neural representation of orientations drifts in WM. This was reflected in a shift of the
87 reconstructed orientation towards the end of the maintenance period that predicted behaviour.

88 **Methods**

89 **Participants**

90 Twenty-six healthy adults (17 female, mean age 25.8 years, range 20-42 years) were included
91 in all analyses. Four additional participants were excluded during preprocessing due to
92 excessive eye-movements (more than 30% of trials contaminated). Participants received
93 monetary compensation (£10 an hour) for participation and gave written informed consent. The
94 experiment was approved by the Central University Research Ethics Committee of the
95 University of Oxford.

96 **Apparatus and stimuli**

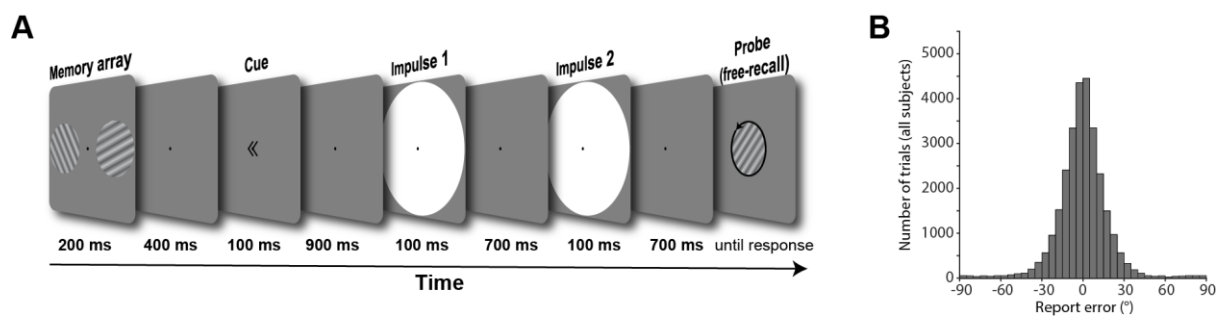
97 The experimental stimuli were generated and controlled by Psychtoolbox (19), a freely
98 available Matlab extension. Visual stimuli were presented on a 23-inch (58.42 cm) screen
99 running at 100 Hz and a resolution of 1,920 by 1,080. Viewing distance was set at 64 cm. A
100 Microsoft Xbox 360 controller was used for response input by the participants.

101 A grey background (RGB = 128, 128, 128; 20.5 cd/m²) was maintained throughout the
102 experiment. A black fixation dot with a white outline (0.242°) was presented in the centre of
103 the screen throughout all trials. Memory items and the probe were sine-wave gratings presented
104 at 20% contrast, with a diameter of 8.51° and spatial frequency of 0.65 cycles per degree, with
105 randomised phase within and across trials. Memory items were presented at 6.08° eccentricity.
106 The rotation of memory items and probe were randomized individually for each trial. The
107 impulse stimulus was a single white circle, with a diameter of 20.67°, presented at the centre
108 of the screen. The retro-cue was two arrowheads pointing right (>>) or left (<<), and was 1.58°
109 wide. A coloured circle (3.4°) was used for feedback. Its colour depended dynamically on the
110 precision of recall, ranging from red (more than 90 degrees error) to green (0 degrees error). A
111 pure tone also provided feedback on recall accuracy after each response, ranging from 200 Hz
112 (more than 90 degrees error) to 1,100 Hz (0 degrees error).

113 Procedure

114 Participants participated in a free-recall, retro-cue visual WM task. Each trial began with the
115 fixation dot. After 1,000 ms the memory array was presented for 200 ms. After a 400 ms delay,
116 the retro-cue was presented for 100 ms, indicating which of the previously two items would be
117 tested, rendering the other item irrelevant. The first impulse stimulus was presented for 100
118 ms, 900 ms after the offset of the retro-cue. After a delay of 700 ms, the second impulse
119 stimulus was presented for 100 ms. After another delay of 700 ms the probe was presented.
120 Participants used the left joystick on the controller with the left thumb to rotate the orientation
121 of the probe until it best reflected the memorized orientation, and confirmed their answer by
122 pressing the “x” button on the controller with the right thumb. Note that one complete rotation
123 of the joystick corresponded to 0.58 of a rotation of the probe. In conjunction with the fact that
124 the probe was randomly orientated on each trial, it was impossible for participants to plan the
125 rotation beforehand or memorize the direction of the joystick instead of the orientation of the
126 memory item. Accuracy feedback was given immediately after the response where both the
127 coloured circle and tone were presented simultaneously. Each participant completed 1,100
128 trials in total, over a course of approximately 135 minutes, including breaks. See Figure 2A for
129 a trial schematic.

130



131 **Figure 2.** Trial schematic and behavioural results (A) Two randomly orientated grating stimuli
132 were presented laterally. A retro-cue then indicated which of those two would be tested at the
133 end of the trial. Two impulses (white circles) were serially presented in the subsequent delay
134 period. At the end of the trial a randomly oriented probe grating was presented in the centre of
135 the screen, and participants were instructed to rotate this probe until it reflected the cued
136 orientation. (B) Report errors of all trials across all subjects.

137 **EEG acquisition**

138 EEG was acquired with 61 Ag/AgCl sintered electrodes (EasyCap, Herrsching, Germany) laid
139 out according to the extended international 10–20 system and recorded at 1,000 Hz using Curry
140 7 software (Compumedics NeuroScan, Charlotte, NC). The anterior midline frontal electrodes
141 (AFz) was used as the ground. Bipolar electrooculography (EOG) was recorded from
142 electrodes placed above and below the right eye and the temples. The impedances were kept
143 below 5 k Ω . The EEG was referenced to the right mastoid during acquisition.

144 **EEG preprocessing**

145 Offline, the EEG signal was re-referenced to the average of both mastoids, down-sampled to
146 500 Hz, and bandpass filtered (0.1 Hz high-pass and 40 Hz low-pass) using EEGLAB (20).
147 The continuous data was epoched relative to the memory array onset (-500 ms to 3,600 ms)
148 before independent component analysis (21) was applied. Components related to eye-blinks
149 were subsequently removed. The data was then epoched relative to memory array onset and
150 the two impulse onsets (0 ms to 400 ms), and trials were individually inspected. Trials with
151 saccadic eye movements, visually identified from the electrooculography, and trials with non-
152 archetypical artefacts, visually identified from the EEG, in the memory array epoch and in
153 either impulse epoch were removed from all subsequent analyses. Furthermore, trials where
154 the report error was 3 circular standard deviations from the participant's mean response error
155 were also excluded from EEG analyses to remove trials that likely represent complete guesses
156 (22). This lead to the removal of $M = 2.3\%$ ($SD = 1.2\%$) trials due to inaccurate report trials,
157 in addition to the $M = 3.52\%$ ($SD = 4.21\%$) and $M = 5\%$ ($SD = 5.2\%$) of trials removed due to
158 eye-movements and non-archetypical EEG artefacts from the memory array and impulse
159 epochs, respectively.

160 While MVPA on electrophysiological data is usually performed on each time-point separately,
161 taking advantage of the highly dynamic waveform of evoked responses in EEG by pooling
162 information multivariately over electrodes as well as time can improve decoding accuracy, at
163 the expense of temporal resolution (23,24). Since the previously reported WM-dependent
164 impulse response reflects the interaction of the WM state at the time of stimulation and does
165 not reflect continuous delay activity, we treat the impulse responses as discrete events in the
166 current study. Thus, the whole time window of interest relative to impulse onsets (100 to 400
167 ms) from the 17 posterior channels was included in the analysis. The time window was based
168 on previous, time-resolved findings, which showed that the WM-dependent neural response

169 from a 100 ms impulse (as used in the current study) is largely confined to this window (18).
170 In the current study, instead of decoding at each time-point separately, information was pooled
171 across the whole time-window. The mean activity level within each time window of each
172 channel was first removed, thus normalizing the voltage fluctuations and isolating the dynamic,
173 impulse-evoked neural signal from more stable brain states. The time-window was then down-
174 sampled by taking the average every 10 ms, thus resulting in 50 values per channel, each of
175 which was treated as a separate dimension in the subsequent multivariate analysis (850 in total).
176 This data format was used on all subsequent MVPA analyses, unless explicitly mentioned
177 otherwise. The same approach over the same time window of interest was used in our previous
178 study (25).

179 **Orientation reconstruction**

180 We computed the mahalanobis distances as a function of orientation difference to reconstruct
181 grating orientations (18). The following procedure was performed separately for items that
182 were presented on the left and right side. Since the grating orientations were determined
183 randomly on a trial-by-trial basis and the resulting orientation distribution across trials was
184 unbalanced, we used a k-fold procedure with subsampling to ensure unbiased decoding. Trials
185 were first assigned the closest of 16 orientations (variable, see below) which were then
186 randomly split into 8 folds using stratified sampling. Using cross-validation, the train trials in
187 7 folds were used to compute the covariance matrix using a shrinkage estimator (26). The
188 number of trials of each orientation bin were equalized by randomly subsampling the minimum
189 number of trials in any bin. The subsampled trials of each angle bin were then averaged. To
190 pool information across similar orientations, the average bins were convolved with a half
191 cosine basis set raised to the 15th power (27). The mahalanobis distances between each trial of
192 the left-out test fold and the averaged and basis-weighted angle bins were computed and mean-
193 centred across the 16 distances to normalize. This was repeated for all test and train fold
194 combinations. To get reliable estimates, the above procedure was repeated 50 times (random
195 folds and subsamples each time), separately for eight orientation spaces (0° to 168.75° , 2.8125°
196 to 171.5625° , 5.625° to 174.375° , 8.4375° to 177.1875° , each in steps of 11.25°). For each trial
197 we thus obtained 200 samples for each of the 16 mahalanobis distances. The distances were
198 averaged across the samples of each trial and ordered as a function of orientation difference.
199 The resulting “tuning curve” was summarized into a single value (i.e., “decoding accuracy”)
200 by computing the cosine vector mean of the tuning curve (18), where a positive value suggests

201 a higher pattern similarity between similar orientations than between dissimilar orientations.
202 The approach was the same for the reanalysis of Wolff et al. (2015).

203 We also repeated the above analysis iteratively for a subset of electrodes in a searchlight
204 analysis across all 61 electrodes. In each iteration, the “current” as well as the closest two
205 neighbouring electrodes were included in the analysis (similar as in: Ede, Chekroud, Stokes, &
206 Nobre, 2019). The freely available MATLAB extension fieldtrip (29) was used to visualise the
207 decoding topographies. Note that the topographies were flipped, such that the left represents
208 the ipsilateral and the right the contralateral side relative to stimulus presentation side.

209 **Orientation code generalization**

210 To test cross-generalization between impulses, instead of training and testing within the same
211 time-window, the train folds were taken from impulse 1, and the test fold from impulse 2, and
212 vice versa. The analysis was otherwise exactly as described above.

213 To test cross-generalization between presented locations, the classifier was similarly trained on
214 trials where the item was presented on the left, and tested on the right, and vice versa. Since
215 left and right trials were independent trial sets, cross-validation does not apply. However, to
216 ensure a balanced training set, the number of trials of each orientation bin were nevertheless
217 equalized by subsampling (as described above), and this approach was repeated 50 times.

218 The cross-generalization of the orientation code between impulse onsets in Wolff et al. (2015)
219 was tested with the same analyses as the location cross-generalization described in the
220 paragraph above: The classifier was trained on the early onset condition, and tested on the late-
221 onset condition, and vice versa, while making sure that the training set is balanced through
222 random subsampling.

223 **Impulse/time and location decoding**

224 To decode the difference of the evoked neural responses between impulses, we used a leave-
225 one-out approach. The mahalanobis distances between the signals from a single trial from both
226 impulse epochs and the average signal of all other trials of each impulse epoch were computed.
227 The covariance matrix was computed by concatenating the trials of each impulse (excluding
228 the left-out trial). The average difference of same impulse distances were subsequently
229 subtracted from different impulse distances, such that a positive distance difference indicates
230 more similarity between same than different impulses. To convert the distance difference into
231 trial wise decoding accuracy, positive distance difference were simply converted into “hits” (1)

232 and negative into “misses” (0). The percentage of correctly classified impulses were
233 subsequently compared to chance performance (50%).

234 The presentation side and impulse onset (in Wolff et al., 2015) was decoded using 8-fold cross-
235 validation, where the distance difference between different and same location/onset was
236 computed for each trial, which were then converted to “hits” and “misses”.

237 **Visualization of the spatial, temporal, and orientation code**

238 To explore and visualize the relationship between the location or impulse/time code and the
239 orientation code in state space (see Fig. 1A for different predictions), we used classical
240 multidimensional scaling (MDS) of the mahalanobis distances between the average signal of
241 trials belonging to one of four orientation bins (0° to 45° , 45° to 90° , 90° to 135° , 135° to 180°)
242 and location (left/right) or time (impulse 1/impulse2).

243 For the visualization of the code across impulse/time, distances were computed separately for
244 left and right trials, before taking the average. Within each orientation bin, the data of half of
245 the trials were taken from impulse 1, and the data of the other half from impulse 2 (determined
246 randomly). The number of trials within each orientation of each impulse were equalized
247 through random subsampling before averaging. The mahalanobis distances between both
248 orientation and impulses were then computed using the covariance matrix estimated from all
249 trials of both impulses. This was repeated 50 times (for each side), randomly subsampling and
250 splitting trials between impulses each time and then taking the average across all iterations.

251 For the visualization of the code across space, the data of each trial were first averaged across
252 impulses. The number of trials of orientation bins (same as above) of each location were
253 equalized through random subsampling. The mahalanobis distances of the average of each bin
254 within each location condition were computed using covariance estimated from all left and
255 right trials. This was repeated 50 times, before taking the average across all iterations.

256 For the code across impulse onset/time visualization of the data from Wolff et al. (2015), the
257 same procedure as in the paragraph above was used, but instead of visualizing the stimulus
258 code between locations, it was visualized between impulse onsets (-30 ms, +30 ms).

259 **Relationship between behaviour and the neural representation of the WM item**

260 We were interested if imprecise reports that are clockwise (CW) or counter-clockwise (CCW)
261 relative to the actual orientation are accompanied by a corresponding shift of the neural

262 representation in WM (see Fig. 1B for model schematics). We used two approaches to test for
263 such a shift (Fig. 5A & 6A).

264 First, the trial-wise pattern similarities as a function of orientation differences (as obtained from
265 the orientation-reconstruction approach described above) were averaged separately for all CW
266 and CCW responses (Fig. 5A). Note that CW and CCW responses were defined relative to the
267 median response error within each orientation bin. This ensures a balanced proportion of all
268 orientations in CW and CCW trials, which is necessary to obtain meaningful orientation
269 reconstructions. It furthermore removes the report bias away from cardinal angles in the current
270 experiment (Suppl. fig. 1), similar to previous reports of orientation response biases (30), and
271 thus isolates random from systematic report errors.

272 We used another approach that exaggerates the potential difference between CW and CCW
273 trials and thus might be more sensitive to detect a shift. The data was first divided into CW and
274 CCW trials using the same within orientation bin approach as described above. The classifier
275 was then trained on CW trials, and tested on CCW trials, and vice versa (Fig. 6A). The
276 orientation bins in the training set were balanced through random subsampling, and the
277 procedure was repeated 50 times. Given an actual shift in the neural representation, the shift
278 magnitude of the resulting orientation reconstruction of this method should be doubled, since
279 both the testing data and the training data (the reference point) are shifted, but in opposite
280 directions.

281 To improve orientation reconstruction from the impulse epochs, the classifier was trained on
282 the averaged trials of both impulses, but tested separately on each impulse epoch individually.
283 While training on both impulses improved orientation reconstruction, in particular for the
284 second approach where only half of the trials are used for training, the shifts in orientation
285 representations as a function of CW/CCW reports are qualitatively the same when training and
286 testing within each impulse epoch separately (Fig. 5, 6, & Suppl. fig. 3).

287 **Statistical significance testing**

288 To test for statistical significance of average decoding at the group level, the sign of the data
289 of each participant was randomly flipped with a probability of 50% 100.000 times, and the
290 resulting null-distribution was used to calculate the p value of the null hypothesis (no
291 difference, chance decoding). Note that tests of within condition decoding (within presentation
292 location, impulse/onset) were one-sided, since only positive decoding is plausible in those
293 cases, whereas tests of cross-generalization between conditions were two-sided, since negative

294 decoding is theoretically plausible in those cases. Comparisons of decodability between
295 conditions/items were also two-sided.

296 The possible shift in representation towards the response was quantified and tested for
297 statistical significance at the group level. The circular mean of the shifted average tuning curve
298 (summarized such that a positive shift reflects a shift towards the response) was tested against
299 0. The tuning curve of each subject was flipped left to right with 0.5 probability, such that a
300 subject's positively shifted tuning curve would then be negatively shifted, before computing
301 the circular mean of the resulting tuning curve averaged over all subjects 100,000 times. The
302 resulting null distribution was used to obtain the p-value by calculating the proportion of
303 permuted tuning curves with circular means more positive than the actual group-level circular
304 mean. The test obtained p-value was one-sided, since we expected the shift of the neural
305 representation of the orientation to be towards the response.

306 **Code and data availability**

307 All custom Matlab scripts and data used to generate the main results of this article will be made
308 publicly available upon peer-reviewed publication.

309 **Results**

310 **Item and WM content-specific evoked responses during encoding and maintenance**

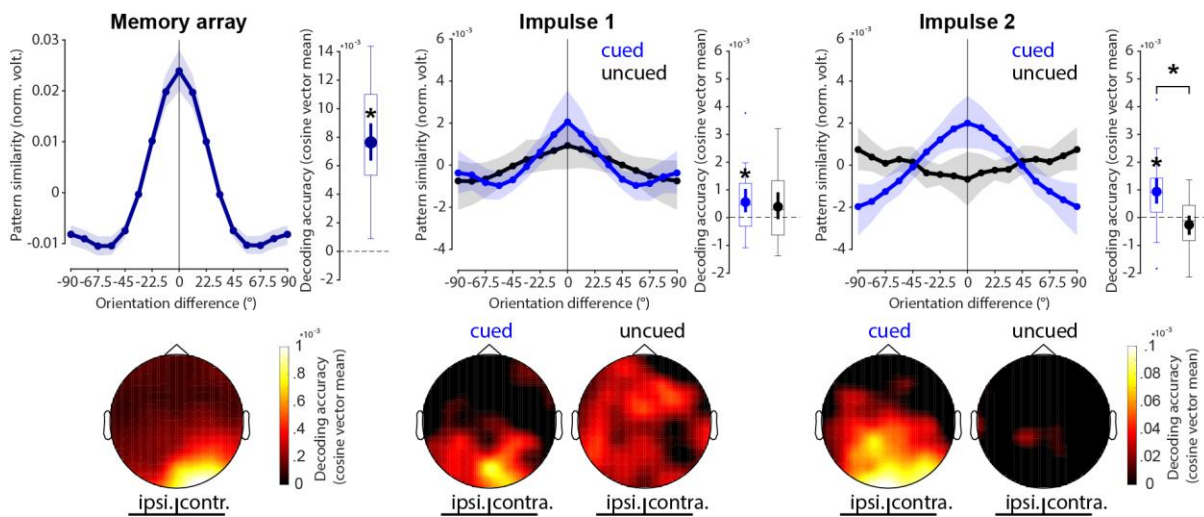
311 The neural response elicited by the memory array presentation contained parametric
312 information about the presentation orientations ($p < 0.001$, one-sided; Fig. 3, left).

313 The first impulse response contained statistically significant information about the cued item
314 ($p = 0.006$, one sided), but not the uncued item, which failed to reach the statistical significance
315 threshold ($p = 0.064$, one-sided). The difference between cued and uncued item decoding was
316 not significant ($p = 0.637$, two-sided; Fig. 3, middle).

317 The decodability of the cued item was also significant at the second impulse response ($p <$
318 0.001 , one-sided), while it was not of the uncued item ($p = 0.917$, one-sided). Notably, the
319 decodability of the cued item was significantly higher than that of the uncued item ($p = 0.001$,
320 two-sided; Fig. 3, right).

321 Overall, these results largely reflect previous findings (18) in that the impulse response reflects
322 relevant information in WM, and that no longer relevant information leave no detectable trace
323 in the WM network.

324 The decoding topographies highlight that most of the decodable signal came from posterior
 325 electrodes during both encoding and maintenance, and is therefore likely generated by the
 326 visual cortex. Notably, while contralateral electrodes showed unsurprisingly higher item
 327 decoding during encoding, this was not the case during maintenance in either impulse response
 328 (Fig. 2C bottom row).
 329



330

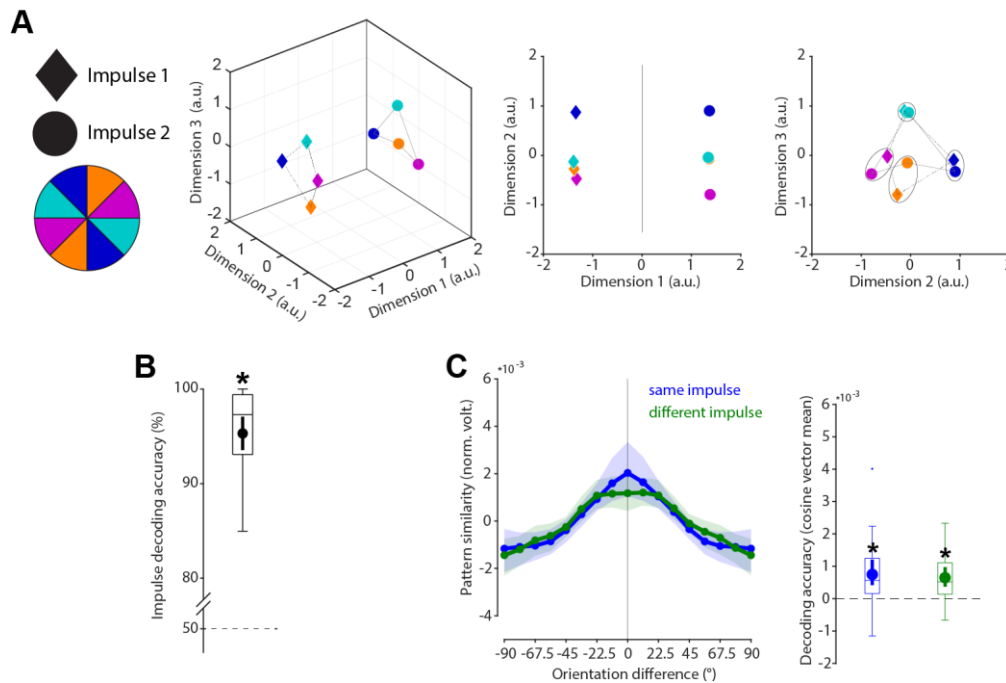
331 **Figure 3.** Decoding results. Top row: Normalized average pattern similarity (mean-centred,
 332 sign-reversed mahalanobis distance) of the evoked neural responses (100 to 400 ms relative to
 333 stimulus onset) as a function of orientation similarity, and decoding accuracy (cosine vector
 334 means of pattern similarities). Error shadings and error bars are 95 % C.I. of the mean. Centre
 335 lines of boxplots indicate the median; box outlines show 25th and 75th percentiles, and
 336 whiskers indicate 1.5x the interquartile range. Extreme values are shown separately (dots).
 337 Asterisks indicate significant decoding accuracies ($p < 0.05$, one-sided) or differences ($p <$
 338 0.05 , two-sided). Bottom row: Decoding topographies of the searchlight analysis.

339 **Stable WM coding scheme in time**

340 The relationship between orientations and impulses/time is visualized in state-space through
 341 MDS (Fig. 4A). While the first dimension clearly differentiates between impulses, the second
 342 and third dimensions code the circular geometry of orientations in both impulses, suggesting
 343 that while the impulse responses are different between impulses, the orientation coding
 344 schemes revealed by the impulse are the same. This is corroborated by significant decoding
 345 accuracy of the impulse ($p < 0.001$, one-sided; Fig. 4B) on the one hand, but also significant

346 cross-generalization of the orientation code between impulses ($p < 0.001$, two-sided), which
 347 was not significantly different from same-impulse orientation decoding ($p = 0.618$, two-sided;
 348 Fig. 4C).

349



350 **Figure 4.** Cross-generalization of coding scheme between impulses. **(A)** Visualization of
 351 orientation and impulse code in state-space. The first dimension discriminates between
 352 impulses. The second and third dimensions code the orientation space in both impulses. **(B)**
 353 Trial-wise accuracy (%) of impulse decoding. **(C)** Orientation decoding within each impulse
 354 (blue) and orientation code cross-generalization between impulses (green). Error shadings and
 355 error bars are 95 % C.I. of the mean. Centre lines of boxplots indicate the median; box outlines
 356 show 25th and 75th percentiles, and whiskers indicate 1.5x the interquartile range. Extreme
 357 values are shown separately (dots). Asterisks indicate significant decoding accuracies or cross-
 358 generalization ($p < 0.05$).

359 It is not possible to conclude whether the difference between impulses is due to a neural
 360 network that changes during the maintenance period over time, due to different stimulation
 361 histories at the time of perturbation (i.e., the first impulse always preceded the second impulse),
 362 or due to different WM operations at each impulse event (e.g. item selection at impulse 1,
 363 response preparation at impulse 2).

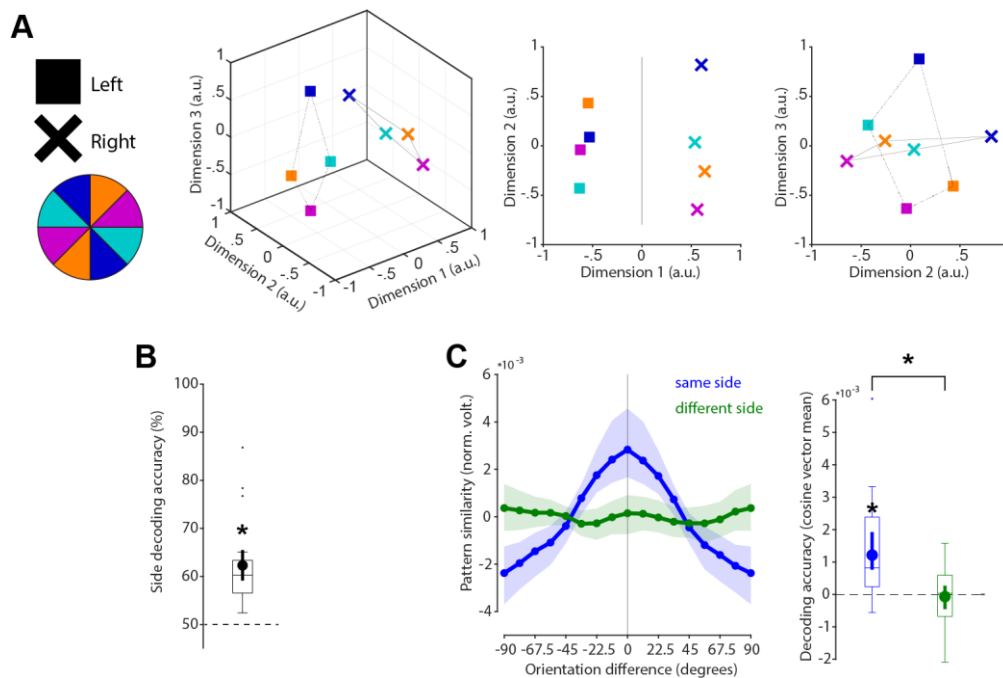
364 To rule out that the difference in impulse response reported above is not only due to difference
365 in stimulation history and changing WM operations, but also due to temporal coding in the
366 WM network, we reanalysed previously published data where a single impulse stimulus was
367 presented either 1,170 or 1,230 ms after the presentation of a single memory item (17). The
368 findings largely replicate the results reported above: State-space visualization of impulse-onset
369 and orientations shows the same circular geometry of the orientations at each impulse onset,
370 while also highlighting a separation of impulse onsets in state-space (Suppl. fig. 2A). Decoding
371 impulse-onset was significantly than from chance ($p = 0.005$, one-sided; Suppl. fig. 2B). Cross-
372 generalization of the orientation code between impulse-onsets was significant ($p < 0.001$, two-
373 sided), and did not significantly differ from decoding the memorized orientation within the
374 same impulse-onset ($p = 0.233$, two-sided; Suppl. fig. 2C).

375 Overall, the results of the current study, as well as the reanalyses of Wolff et al. (2015) provide
376 evidence for a low-dimensional change over time, that can be revealed by perturbing the WM
377 network at different time-points (as predicted in Buonomano & Maass, 2009), while at the
378 same time providing evidence for a temporally stable coding scheme of WM content (3,4).

379 **Specific WM coding scheme in space**

380 As a counterpart to the stable coding scheme in time reported above, we explicitly tested if the
381 coding scheme is location specific (i.e., dependent on the previous presentation location of the
382 cued orientation). State-space visualization of cued item location and orientations shows a clear
383 separation between locations and no overlap in orientation coding between locations (Fig. 5A).
384 The cued location was significantly decodable from the impulse responses ($p < 0.001$, one-
385 sided; Fig. 5B). Cross-generalization of the orientation coding scheme between cued item
386 locations was not significant ($p = 0.716$, two-sided), and significantly lower than same side
387 orientation decoding ($p = 0.002$, two-sided; Fig. 5C). These results reflect previous reports of
388 spatially specific WM codes, even when location is no longer relevant (32).

389



390 **Figure 5.** No cross-generalization of coding scheme between cued item locations during
 391 impulse responses **(A)** Visualization of orientation and item location code in state-space. The
 392 first dimension discriminates between item locations. The first and second dimensions code
 393 the orientation space, separately for WM items previously presented on the left or right side.
 394 **(B)** Trial-wise accuracy (%) of item location decoding. **(C)** Orientation decoding within each
 395 item location (blue) and orientation code cross-generalizing between different item locations
 396 (green). Error shadings and error bars are 95 % C.I. of the mean. Centre lines of boxplots
 397 indicate the median; box outlines show 25th and 75th percentiles, and whiskers indicate 1.5x
 398 the interquartile range. Extreme values are shown separately (dots). Asterisks indicate
 399 significant decoding accuracies and differences ($p < 0.05$).

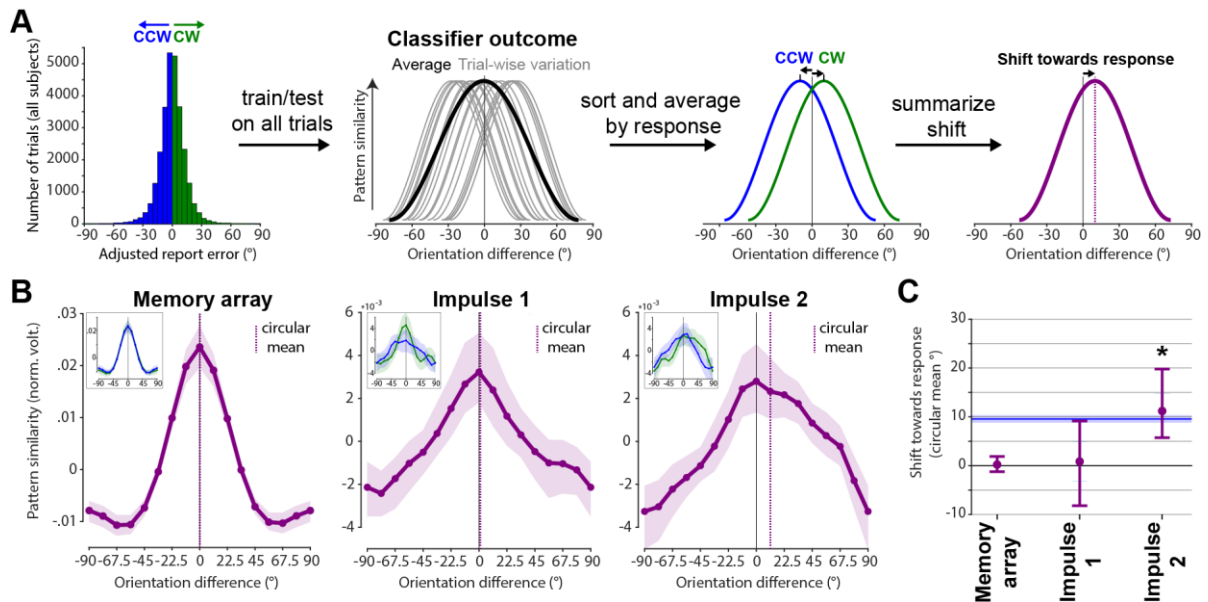
400 **Drifting WM code**

401 The first approach to test for a possible shift of the neural representation towards the response
 402 averaged the trial-wise orientation tuning curves obtained from the cross-validated orientation
 403 reconstruction on all trials (see Methods and Fig. 6A).

404 No significant shift towards the response was evident during encoding/memory array
 405 presentation ($p = 0.394$, one-sided; Fig. 6B & C, left). No evidence for such a shift was found
 406 at impulse 1/early maintenance either ($p = 0.423$, one-sided; Fig. 6B & C, middle). However,

407 the orientation tuning curve was significantly shifted towards the response at impulse 2/late
 408 maintenance ($p < 0.001$, one-sided; Fig. 6B & C, right).

409



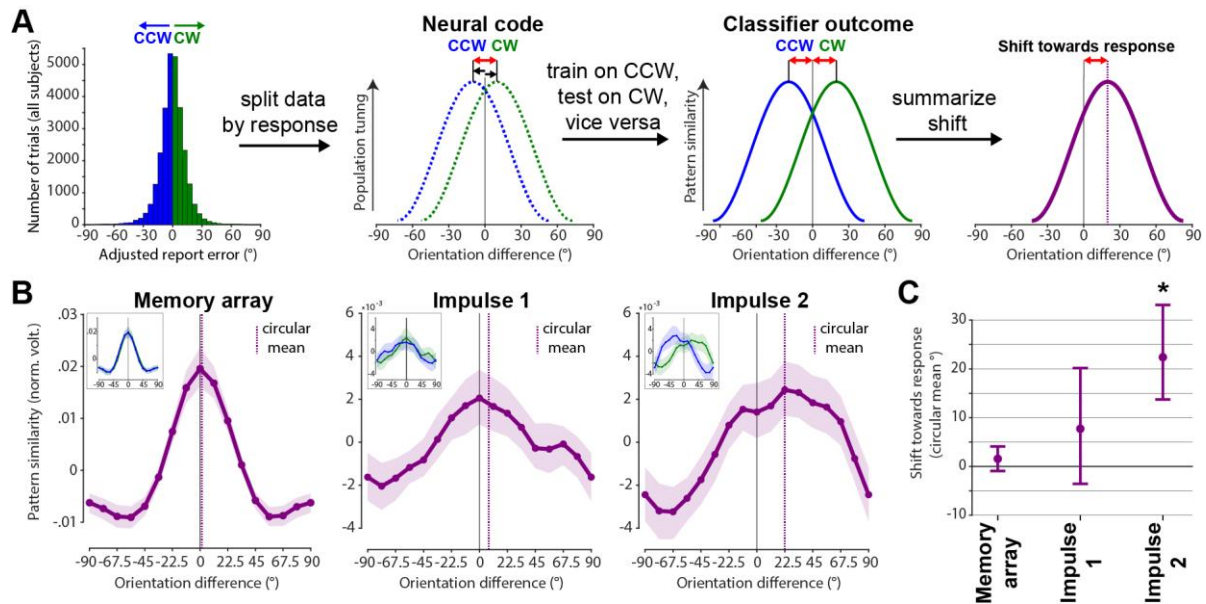
410 **Figure 6.** Response-dependent averaging of trial-wise tuning curves demonstrates drift.
 411 Schematic and results. **(A)** Testing for shift towards response by averaging trial-wise tuning
 412 curves by CCW/CW responses. **(B)** Results of schematised approach in A. Orientation tuning
 413 curves averaged by response such that a right-ward shift reflects a shift towards the response
 414 (purple) at each event. Purple vertical lines show circular means of the tuning curves. Insets
 415 show orientation tuning curves for CCW (blue) and CW (green) responses separately. Error
 416 shadings are 95 % C. I. of the mean. **(C)** Group-level shifts towards the response (circular
 417 mean) of each response-dependent tuning curve. Error-bars are 95 % C. I. of the mean. The
 418 blue line and shading indicates the mean and 95 % C.I. of the absolute, bias-adjusted
 419 behavioural response deviation.

420 The second approach to test for a possible shift of the neural representation towards the
 421 response may be more sensitive since it trains the orientation classifier only on CCW trials,
 422 and tests it on CW trials, and vice versa (see Methods and Fig. 7A), thus increasing any
 423 response related shift by a factor of two.

424 This approach yielded similar results as the previous approach, though the shift magnitudes are
 425 indeed larger. Neither the memory array presentation/encoding, nor impulse 1/early
 426 maintenance showed a significant shift towards the response ($p = 0.121$, $p = 0.095$, respectively,

427 one-sided; Fig. 7, left & middle), while impulse 2/late maintenance did ($p < 0.001$, one-sided;
 428 Fig. 7, right).

429



430 **Figure 7.** Response-dependent training and testing demonstrates drift. Schematic and results.
 431 **(A)** Testing for shift towards response by first splitting the neuroimaging data into CW and
 432 CCW data sets, and training on CW trials and testing on CCW trials, and vice versa. Given an
 433 actual shift, the shift of the resulting orientation reconstruction will be doubled, since training
 434 and testing data are shifted in opposite directions. **(B)** Results of schematised approach in A.
 435 Average orientation tuning curves such that a right-ward shift reflects a shift towards the
 436 response (purple) at each event. Purple vertical lines show circular means of the tuning curves.
 437 Insets show orientation tuning curves for CCW (blue) and CW (green) responses separately.
 438 Error shadings are 95 % C. I. of the mean. **(C)** Group-level shifts towards the response (circular
 439 mean) of each response-dependent tuning curve. Error-bars are 95 % C. I. of the mean.

440 Note the reported results of shifts during impulse presentations were obtained by training the
 441 classifier on both impulses, but testing it on each impulse separately. This was done to improve
 442 power (as explained in Methods). This improved orientation reconstruction particularly for the
 443 latter shift-analysis where the classifier is trained on only half the trials (CW trials only or
 444 CCW trials only). However, the same analyses based on training (and testing) within each
 445 impulse epoch separately yielded qualitatively similar results (no significant shifts at impulse
 446 1 in either approach, significant shifts at impulse 2 in both approaches; Suppl. fig. 3).

447

Discussion

448 In the present study, we investigated the neural dynamics of WM by probing the coding scheme
449 over time, as well as drift in the actual memories. The neural response to impulse stimuli in
450 this non-spatial WM paradigm enabled us to show that the coding scheme of parametric visual
451 feature (i.e., orientation) in WM remained stable during maintenance, reflected in the
452 significant cross-generalization of the orientation decoding between early and late impulses
453 (Fig. 4). However, memories drift within this stable coding scheme, leading to a bias in
454 memories (Figs. 6 and 7).

455 This is consistent with previous reports of a stable subspace for WM maintenance (4,5), and
456 provides evidence for a time-invariant coding scheme for orientations maintained in WM.
457 However, more dynamic schemes have also been reported. For example, during the early
458 transition between encoding and maintenance (33). At the extreme end, some have proposed
459 that WM could be maintained in a dynamical system, where activity evolves along a complex
460 trajectory in neural state space (e.g. 34). Although this complicates readout (discrimination
461 boundaries at one time point do not generalise to other time-points), such coding schemes
462 evolve naturally from recurrent neural networks. Moreover, such dynamics also provide
463 additional information, such as elapsed time. In the current study, we find evidence for a hybrid
464 model (3,4): stable decoding of WM, despite dynamic activity over time.

465 Specifically, while there was no cost of cross-generalizing the orientation code between
466 impulses, there was nevertheless a clear difference in the neural pattern between them,
467 suggesting that a separate dynamic neural pattern codes the passage of time. A reanalysis of
468 the data of a previously published study (17) confirmed these results, suggesting that the low-
469 dimensional dynamics code for time per se (rather than impulse number). The significant
470 decodability of impulse onset shows that the WM network changes during the maintenance
471 even within 60 ms, resulting in distinct neural impulse responses at different time-points
472 providing evidence for a neural time-code. Importantly, the low-dimensional representation of
473 elapsed time is orthogonal to the mnemonic subspace, allowing WM representations to be
474 stable. This hybrid of stable and dynamic representations may emerge from interactions
475 between dynamic recurrent neural networks and stable sensory representations (3).

476 Our index of WM-related neural activity was based on an impulse response approach that we
477 previous developed to measure WM-related changes in the functional state of the system,
478 including ‘activity-silent’ WM states (17,18,35). For example, activity states during encoding

479 could result in a neural trace in the WM network through short-term synaptic plasticity
480 resulting in a stable code for maintenance, whereas the time-dimension could be represented
481 in its gradual fading (31,36–38). The stable WM-content coding scheme could also be
482 achieved by low-level activity states that self-sustain a stable code through recurrent
483 connections, a key feature of attractor models of WM (1,39), while dynamic activity patterns
484 are coded in an orthogonal subspace that represents time. While we did not explicitly consider
485 tonic delay activity, it is nonetheless possible that the impulse responses also reflect non-linear
486 interactions with low-level, persistent activity states that are otherwise difficult to measure with
487 EEG. Therefore, we cannot rule out a contribution of persistent activity in the stable coding
488 scheme observed here.

489 We also found evidence that the orientation code itself drifts along the orientation dimension,
490 predicting recall errors. While there was no bias in the neural orientation representation at either
491 encoding or early maintenance, the second impulse towards the end of the maintenance period
492 revealed a code that was shifted towards the direction of response error. This pattern of results
493 is consistent with the drift account of WM, where neural noise leads to an accumulation of
494 error during maintenance, resulting in a still sharp, but shifted (i.e. slightly wrong) neural
495 representation of the maintained information (1,14). While previous neurophysiological
496 recordings from monkey PFC found evidence for drift for spatial information (15), we could
497 demonstrate a shifting representation that more faithfully represents non-spatial WM content
498 that is unrelated to sustained spatial attention or motor preparation, by using lateralized
499 orientations in the present study.

500 Bump attractors have been proposed as an ideal neural mechanism for the maintenance of
501 continuous representations (i.e. space, orientation, colour), where a specific feature is
502 represented by the persistent activity “bump” of the neural population at the feature’s location
503 along the network’s continuous feature space. Neural noise randomly shifts this bump along
504 the feature dimension, while inhibitory and excitatory connections maintain the same overall
505 level of activity and shape of the neural network (40,41). Random walk along the feature
506 dimension is thus a fundamental property of bump attractors, and has been found to explain
507 neurophysiological findings (15). Typically, this is considered within the framework of
508 persistent working memory, however transient bursts of activity could also follow similar
509 attractor dynamics (42,43). For example, the temporary connectivity changes of the memorized
510 WM item may indeed slowly dissolve and become coarser, periodic activity bursts may keep
511 this to a minimum, by periodically reinstating a sharp representation. However, since this

512 refreshing depends on the read-out of a coarse representation, the resulting representation may
513 be slightly wrong and thus shifted. This interplay between decaying silent WM-states that are
514 readout and refreshed by active WM-states also predicts a drifting WM code, without
515 depending on an unbroken chain of persistent neural activity.

516 Moreover, the representational drift does not necessarily have to be random. Modelling of
517 report errors in a free recall colour WM task suggests that an increase of report errors over time
518 may be due to separable attractor dynamics, with a systematic drift towards stable colour
519 representations, resulting in a clustering of reports around specific colour values, in addition to
520 random drift elicited by neural noise (13). The report bias of oblique orientations seen in the
521 present study could be explained by a similar drift towards specific orientations, which would
522 predict an increase of report bias for longer retention periods. However, clear behavioural
523 evidence for such an increase in systemic report errors of orientations is lacking (10). In the
524 present study we isolated random from systematic errors, both as a methodological necessity,
525 but also to be able to conclude that any observed shift is due to random errors. Thus, while a
526 systematic drift towards specific orientations might be possible, the shift in representation
527 reported here is unrelated to it.

528 Our results suggest that maintenance in WM is dynamic, although the fundamental coding
529 scheme remains stable over time. Low-dimensional dynamics could provide a valuable readout
530 of elapsed time, whilst allowing for a time-general readout scheme for the WM content. We
531 also show that drift within this stable coding scheme could explain loss of memory precision
532 over time.

533

Acknowledgments

534 This research was in part funded by a James S. McDonnell Foundation Scholar Award
535 (220020405) to MGS, and by the NIHR Oxford Health Biomedical Research Centre. The
536 Wellcome Centre for Integrative Neuroimaging is supported by core funding from the
537 Wellcome Trust (203139/Z/16/Z). The views expressed are those of the authors and not
538 necessarily those of the National Health Service, the National Institute for Health Research or
539 the Department of Health. EGA is in part funded by an Open Research Area grant (464.18.114).
540 We would like to thank N.E. Myers and D. Trübtschek for helpful comments.

541

References

- 542 1. Compte A, Brunel N, Goldman-Rakic PS, Wang X-J. Synaptic Mechanisms and Network
543 Dynamics Underlying Spatial Working Memory in a Cortical Network Model. *Cereb*
544 *Cortex*. 2000 Sep 1;10(9):910–23.
- 545 2. Wang X-J. Synaptic reverberation underlying mnemonic persistent activity. *Trends in*
546 *Neurosciences*. 2001 Aug 1;24(8):455–63.
- 547 3. Bouchacourt F, Buschman TJ. A Flexible Model of Working Memory. *Neuron*. 2019 Jul
548 3;103(1):147-160.e8.
- 549 4. Murray JD, Bernacchia A, Roy NA, Constantinidis C, Romo R, Wang X-J. Stable
550 population coding for working memory coexists with heterogeneous neural dynamics in
551 prefrontal cortex. *PNAS*. 2017 Jan 10;114(2):394–9.
- 552 5. Spaak E, Watanabe K, Funahashi S, Stokes MG. Stable and Dynamic Coding for Working
553 Memory in Primate Prefrontal Cortex. *J Neurosci*. 2017 Jul 5;37(27):6503–16.
- 554 6. Cueva CJ, Marcos E, Saez A, Genovesio A, Jazayeri M, Romo R, et al. Low dimensional
555 dynamics for working memory and time encoding. *bioRxiv*. 2019 Jan 31;504936.
- 556 7. Barak O, Sussillo D, Romo R, Tsodyks M, Abbott LF. From fixed points to chaos: Three
557 models of delayed discrimination. *Progress in Neurobiology*. 2013 Apr 1;103:214–22.
- 558 8. Romo R, Brody CD, Hernández A, Lemus L. Neuronal correlates of parametric working
559 memory in the prefrontal cortex. *Nature*. 1999 Jun;399(6735):470.
- 560 9. Druckmann S, Chklovskii DB. Neuronal Circuits Underlying Persistent Representations
561 Despite Time Varying Activity. *Current Biology*. 2012 Nov 20;22(22):2095–103.
- 562 10. Rademaker RL, Park YE, Sack AT, Tong F. Evidence of gradual loss of precision for
563 simple features and complex objects in visual working memory. *Journal of Experimental*
564 *Psychology: Human Perception and Performance*. 2018;44(6):925–40.
- 565 11. Barrouillet P, Camos V. Developmental Increase in Working Memory Span: Resource
566 Sharing or Temporal Decay? *Journal of Memory and Language*. 2001 Jul 1;45(1):1–20.

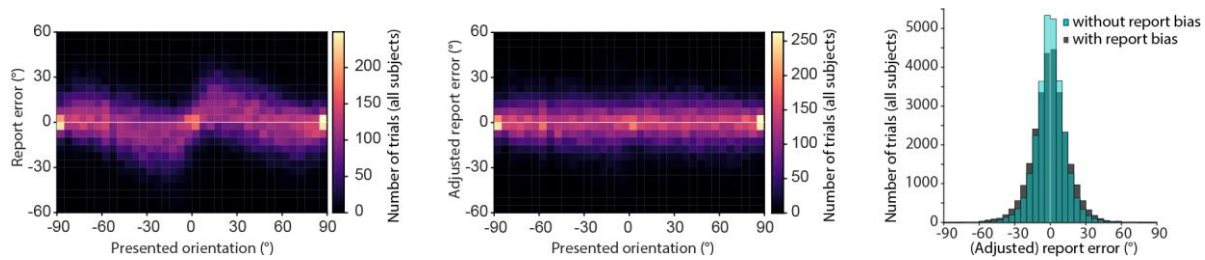
- 567 12. Kinchla RA, Smyzer F. A diffusion model of perceptual memory. *Perception &*
568 *Psychophysics*. 1967 Jun 1;2(6):219–29.
- 569 13. Panichello MF, DePasquale B, Pillow JW, Buschman TJ. Error-correcting dynamics in
570 visual working memory. *Nature Communications*. in press;
- 571 14. Schneegans S, Bays PM. Drift in Neural Population Activity Causes Working Memory
572 to Deteriorate Over Time. *J Neurosci*. 2018 May 23;38(21):4859–69.
- 573 15. Wimmer K, Nykamp DQ, Constantinidis C, Compte A. Bump attractor dynamics in
574 prefrontal cortex explains behavioral precision in spatial working memory. *Nat Neurosci*.
575 2014 Mar;17(3):431–9.
- 576 16. Lim PC, Ward EJ, Vickery TJ, Johnson MR. Not-so-working Memory: Drift in
577 Functional Magnetic Resonance Imaging Pattern Representations during Maintenance
578 Predicts Errors in a Visual Working Memory Task. *Journal of Cognitive Neuroscience*.
579 2019 May 21;1–15.
- 580 17. Wolff MJ, Ding J, Myers NE, Stokes MG. Revealing hidden states in visual working
581 memory using electroencephalography. *Front Syst Neurosci*. 2015 Jul 28;9.
- 582 18. Wolff MJ, Jochim J, Akyürek EG, Stokes MG. Dynamic hidden states underlying
583 working-memory-guided behavior. *Nature Neuroscience*. 2017 Jun;20(6):864–71.
- 584 19. Kleiner M. Visual stimulus timing precision in Psychtoolbox-3: Tests, pitfalls solutions
585 [Internet]. 2010 [cited 2016 Jul 12]. Available from: [http://www.neuroschool-tuebingen-](http://www.neuroschool-tuebingen-ena.de/fileadmin/user_upload/Dokumente/neuroscience/AbstractbookNeNa2010u.pdf)
586 [ena.de/fileadmin/user_upload/Dokumente/neuroscience/AbstractbookNeNa2010u.pdf](http://www.neuroschool-tuebingen-ena.de/fileadmin/user_upload/Dokumente/neuroscience/AbstractbookNeNa2010u.pdf)
- 587 20. Delorme A, Makeig S. EEGLAB: an open source toolbox for analysis of single-trial EEG
588 dynamics including independent component analysis. *Journal of Neuroscience Methods*.
589 2004 Mar;134(1):9–21.
- 590 21. Hyvarinen A. Fast and robust fixed-point algorithms for independent component analysis.
591 *IEEE Transactions on Neural Networks*. 1999 May;10(3):626–34.
- 592 22. Fritsche M, Mostert P, de Lange FP. Opposite Effects of Recent History on Perception
593 and Decision. *Current Biology*. 2017 Feb 20;27(4):590–5.

- 594 23. Grootswagers T, Wardle SG, Carlson TA. Decoding Dynamic Brain Patterns from
595 Evoked Responses: A Tutorial on Multivariate Pattern Analysis Applied to Time Series
596 Neuroimaging Data. *J Cogn Neurosci*. 2017 Apr;29(4):677–97.
- 597 24. Nemrodov D, Niemeier M, Patel A, Nestor A. The Neural Dynamics of Facial Identity
598 Processing: Insights from EEG-Based Pattern Analysis and Image Reconstruction.
599 *eNeuro*. 2018 Jan 1;5(1):ENEURO.0358-17.2018.
- 600 25. Wolff MJ, Kandemir G, Stokes MG, Akyurek EG. Impulse responses reveal unimodal
601 and bimodal access to visual and auditory working memory. *bioRxiv*. 2019 Apr
602 30;623835.
- 603 26. Ledoit O, Wolf M. Honey, I shrunk the sample covariance matrix. *The Journal of Portfolio
604 Management*. 2004;30(4):110–9.
- 605 27. Myers NE, Rohenkohl G, Wyart V, Woolrich MW, Nobre AC, Stokes MG. Testing
606 sensory evidence against mnemonic templates. *eLife*. 2015 Dec 14;4:e09000.
- 607 28. Ede F van, Chekroud SR, Stokes MG, Nobre AC. Concurrent visual and motor selection
608 during visual working memory guided action. *Nature Neuroscience*. 2019 Mar;22(3):477.
- 609 29. Oostenveld R, Fries P, Maris E, Schoffelen J-M, Oostenveld R, Fries P, et al. FieldTrip:
610 Open Source Software for Advanced Analysis of MEG, EEG, and Invasive
611 Electrophysiological Data, FieldTrip: Open Source Software for Advanced Analysis of
612 MEG, EEG, and Invasive Electrophysiological Data. *Computational Intelligence and
613 Neuroscience, Computational Intelligence and Neuroscience*. 2010 Dec 23;2011,
614 2011:e156869.
- 615 30. Pratte MS, Park YE, Rademaker RL, Tong F. Accounting for stimulus-specific variation
616 in precision reveals a discrete capacity limit in visual working memory. *Journal of
617 Experimental Psychology: Human Perception and Performance*. 2017;43(1):6–17.
- 618 31. Buonomano DV, Maass W. State-dependent computations: spatiotemporal processing in
619 cortical networks. *Nature Reviews Neuroscience*. 2009 Feb;10(2):113–25.
- 620 32. Pratte MS, Tong F. Spatial specificity of working memory representations in the early
621 visual cortex. *Journal of Vision*. 2014 Mar 1;14(3):22–22.

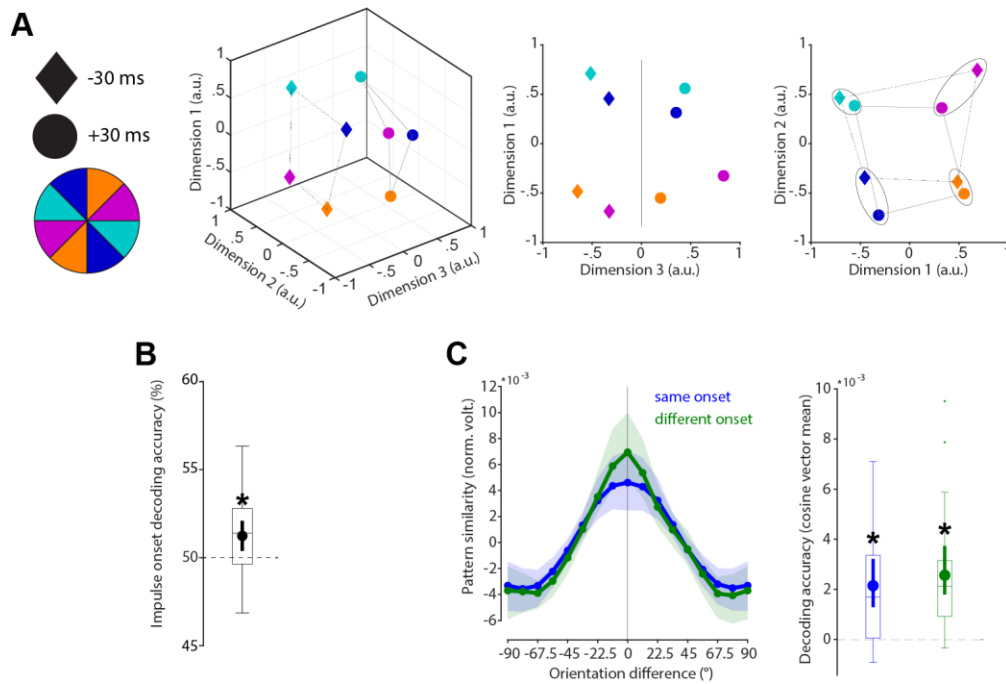
- 622 33. Wasmuht DF, Spaak E, Buschman TJ, Miller EK, Stokes MG. Intrinsic neuronal
623 dynamics predict distinct functional roles during working memory. *Nature*
624 *Communications*. 2018 Aug 29;9(1):3499.
- 625 34. Maass W, Natschläger T, Markram H. Real-time computing without stable states: a new
626 framework for neural computation based on perturbations. *Neural Comput*. 2002
627 Nov;14(11):2531–60.
- 628 35. Stokes MG. ‘Activity-silent’ working memory in prefrontal cortex: a dynamic coding
629 framework. *Trends in Cognitive Sciences*. 2015 Jul;19(7):394–405.
- 630 36. Nikolić D, Häusler, Stefan, Singer, Wolf, Maass, Wolfgang. Temporal dynamics of
631 information content carried by neurons in the primary visual cortex. *Advances in Neural*
632 *Information Processing Systems*. 2007;19:1041–1048.
- 633 37. Nikolić D, Häusler S, Singer W, Maass W. Distributed Fading Memory for Stimulus
634 Properties in the Primary Visual Cortex. *PLOS Biol*. 2009 Dec 22;7(12):e1000260.
- 635 38. Zucker RS, Regehr WG. Short-Term Synaptic Plasticity. *Annu Rev Physiol*. 2002 Mar
636 1;64(1):355–405.
- 637 39. Chaudhuri R, Fiete I. Computational principles of memory. *Nature Neuroscience*. 2016
638 Mar;19(3):394–403.
- 639 40. Amari S. Dynamics of pattern formation in lateral-inhibition type neural fields. *Biol*
640 *Cybern*. 1977 Jun 1;27(2):77–87.
- 641 41. Brody CD, Romo R, Kepecs A. Basic mechanisms for graded persistent activity: discrete
642 attractors, continuous attractors, and dynamic representations. *Current Opinion in*
643 *Neurobiology*. 2003 Apr 1;13(2):204–11.
- 644 42. Lundqvist M, Herman P, Lansner A. Theta and Gamma Power Increases and Alpha/Beta
645 Power Decreases with Memory Load in an Attractor Network Model. *Journal of*
646 *Cognitive Neuroscience*. 2011 Mar 31;23(10):3008–20.
- 647 43. Mongillo G, Barak O, Tsodyks M. Synaptic Theory of Working Memory. *Science*. 2008
648 Mar 14;319(5869):1543–6.

649

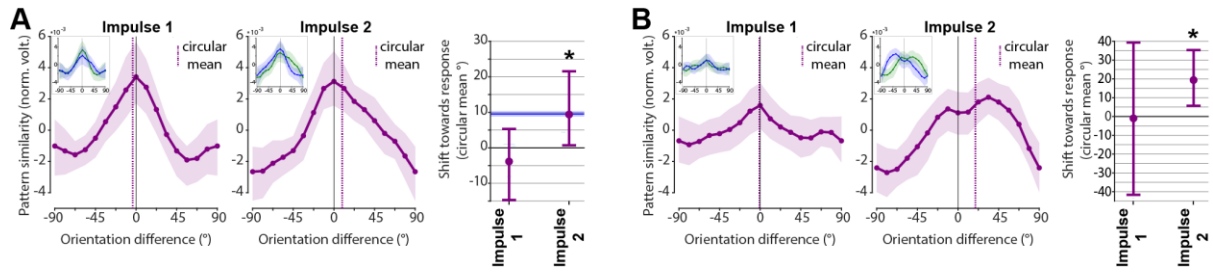
Appendix



650 **Supplemental figure 1.** Report-bias of orientations. Participants showed a bias, exaggerating
651 the tilt of oblique orientations, manifesting itself as a repulsion from the cardinal axes (0 and
652 90 degrees; *left*), similar to previous reports (30). To ensure an unbiased estimate of a possible
653 shift in our analysis, and to isolate random from systematic errors, the report bias was removed
654 by subtracting the median error within 11.25 degree orientation bins (*middle*). By removing
655 orientation-specific error, the resulting error distribution is narrower (*right*). Clockwise and
656 counter-clockwise reports were defined as positive and negative reports relative to this
657 “adjusted”, unbiased, report error.



658 **Supplemental figure 2.** Cross-generalization of coding scheme between impulse onsets in
659 reanalyses of Wolff et al. (2015). **(A)** Visualization of orientation and impulse-onset code in
660 state-space. The third dimension discriminates between impulse-onsets. The first and second
661 dimensions code the orientation space in both impulses. **(B)** Trial-wise accuracy (%) of
662 impulse-onset decoding. **(C)** Orientation decoding within each impulse-onset (blue) and
663 orientation code cross-generalizing between impulse-onsets (green). Error shadings and error
664 bars are 95 % C.I. of the mean. Centre lines of boxplots indicate the median; box outlines show
665 25th and 75th percentiles, and whiskers indicate 1.5x the interquartile range. Extreme values
666 are shown separately (dots). Asterisks indicate significant decoding accuracies or cross-
667 generalization ($p < 0.05$).



668

669

670

671

672

673

Supplemental figure 3. Within impulse training and testing to estimate drift. **(A)** Response-dependent averaging of trial-wise tuning curves (Fig. 6A). Shift towards response: Impulse 1: $p = 0.807$; Impulse 2: $p = 0.029$, one-sided. **(B)** Response-dependent training and testing (Fig. 7A). Shift towards response: Impulse 1: $p = 0.524$; Impulse 2: $p = 0.007$, one-sided. Same convention as Fig. 6B-C and Fig. 7B-C.