

1 **Current and future goals are represented in opposite patterns in object-**
2 **selective cortex**

3 Anouk M. van Loon^{1,2*}, Katya Olmos Solis^{1*}, Johannes J. Fahrenfort^{1,3§} & Christian N. L. Olivers^{1,2§}

4 ¹Department of Experimental and Applied Psychology, Vrije Universiteit Amsterdam,

5 ²Institute of Brain and Behavior Amsterdam, Vrije Universiteit Amsterdam, Amsterdam, The
6 Netherlands

7 ³Department of Brain and Cognition, University of Amsterdam, The Netherlands

8 Address correspondence to:

9 Christian N. L. Olivers

10 Department of Experimental and Applied Psychology,

11 Vrije Universiteit Amsterdam,

12 Van der Boechorststraat 1,

13 1081 BT, Amsterdam,

14 The Netherlands.

15 Email: c.n.l.olivers@vu.nl

16 *shared first author

17 §shared senior author

18 Acknowledgements:

19 We thank Assaf Harel for providing the pictures used as stimuli. This research was supported by the
20 European Research Council (ERC) under grant agreement no. ERC-CoG-2013-615423 awarded to
21 CNLO.

22

23 Conflict of Interest:

24 None declared.

25

26

27

1 **Abstract**

2 Adaptive behavior requires the separation of current from future goals in working memory. We
3 used fMRI of object-selective cortex to determine the representational (dis)similarities of memory
4 representations serving current and prospective perceptual tasks. Participants remembered an
5 object drawn from three possible categories as the target for one of two consecutive visual search
6 tasks. A cue indicated whether the target object should be looked for first (currently relevant),
7 second (prospectively relevant), or if it could be forgotten (irrelevant). Prior to the first search,
8 representations of current, prospective and irrelevant objects were similar, with strongest
9 decoding for current representations compared to prospective (Experiment 1) and irrelevant
10 (Experiment 2). Remarkably, during the first search, prospective representations could also be
11 decoded, but revealed anti-correlated voxel patterns compared to currently relevant
12 representations of the same category. We propose that the brain separates current from
13 prospective memories within the same neuronal ensembles through opposite representational
14 patterns.

15

16 **Introduction**

17 Adaptive human behavior requires the representation of both imminent and future goals in
18 response to changing task requirements. Little is known about how the brain distinguishes between
19 information that is currently relevant and information that is only prospectively relevant.

20 While working memory is thought to be pivotal to the active maintenance of
21 representations for current task goals, representations serving prospective tasks should be
22 shielded from affecting currently relevant input and output, and vice versa. Studies using reaction
23 time and eye movement measures have indeed shown that currently and prospectively relevant
24 representations differentially bias processing of perceptual input (e.g., Carlisle & Woodman, 2011;
25 Downing & Dodds, 2003; Houtkamp & Roelfsema, 2006; Mallett & Lewis-Peacock, 2018; Olivers &
26 Eimer, 2011; van Loon, Olmos Solis, & Olivers, 2017). Furthermore, studies using multi-variate
27 pattern analyses (MVPA) of functional magnetic resonance imaging (fMRI) and
28 electroencephalography (EEG) data have shown that while representations required for an
29 upcoming memory test can be readily decoded, the evidence for items required for a prospective
30 task temporarily drops to baseline levels until they become relevant again (LaRocque, Lewis-
31 Peacock, Drysdale, Oberauer, & Postle, 2013; LaRocque, Riggall, Emrich, & Postle, 2017; Lewis-
32 Peacock & Postle, 2012). These and other findings have led to the hypothesis that items in working

1 memory may adopt different states or representational formats (Barak & Tsodyks, 2014;
2 D'Esposito & Postle, 2015; LaRocque, Lewis-Peacock, & Postle, 2014; Olivers, Peters, Houtkamp, &
3 Roelfsema, 2011; Stokes, 2015). While currently relevant items are represented through patterns of
4 firing across populations of neurons, prospectively relevant items may be stored in what has been
5 referred to as an “activity-silent”, or “hidden” state. One way in which such a state can be achieved
6 is through short-term potentiation of synaptic connectivity in the neuronal population, as induced
7 by the initial firing activity during encoding and active storage within that same population
8 (Erickson, Maramba, & Lisman, 2010; Mongillo, Barak, & Tsodyks, 2008; Sugase-Miyamoto, Liu,
9 Wiener, Optican, & Richmond, 2008). Another way is through changes in the membrane potentials
10 of the previously firing neurons (e.g., Stokes, 2015). We will collectively refer to these options as
11 changes in the responsivity (versus the activity) of a neuronal ensemble.

12 Such latent changes in responsivity are by definition difficult to test through activity-based
13 measures. One prediction is that prospective memories re-emerge in activity-based dependent
14 measures when *unrelated* activity is sent through the network and interacts with the pattern of
15 changed responsivity that reflects the activity-silent memory. This is indeed what Rose and
16 colleagues (2016) recently reported. They found that prospective memory representations which
17 could initially no longer be decoded during a working memory delay period could successfully be
18 reconstructed after applying a brief burst of transcranial magnetic stimulation (TMS, Rose et al.,
19 2016). Likewise, Wolff and colleagues recently reported enhanced decoding of a memorized
20 oriented grating shortly after observers were presented with a visual pattern that was neutral with
21 respect to the memorized orientation (Wolff, Ding, Myers, & Stokes, 2015; Wolff, Jochim, Akyurek, &
22 Stokes, 2017). However, although these studies show that there is information present on
23 prospectively stored memories, it is as yet unclear what the representational format of such
24 prospective memories is, and how they relate to currently relevant memories.

25 A priori there appear to be a number of hypotheses. First, the standard synaptic
26 potentiation mechanism predicts that the altered pattern of responsivity directly follows the
27 pattern of activity during encoding of the item, thus predicting a high degree of similarity between
28 the active and the silent representation when revived. A second possibility is that it is unnecessary
29 to assume activity-silent representations at all, as has recently been argued by Schneegans & Bays
30 (2017). Instead, they argued for a single maintenance mechanism in which differently prioritized
31 items in memory are stored through similar patterns of firing activity, with the only difference
32 being the degree of activation. Their model simulations provide a proof of concept that the revival
33 of a memory can be explained by selectively boosting the still present, but lowered activity, rather

1 than by the reconstruction from hidden states of responsivity. Also under this scenario the same
2 pattern of activation should emerge for current and prospective memories, except for a difference
3 in strength. The third possibility is that prospectively relevant items are stored in an altogether
4 different pattern compared to actively maintained items – that is, they may be transformed within
5 the same population, or stored in different populations, whether through changed activity or
6 responsivity. This was recently proposed by Christophel and colleagues (2018), who found
7 currently relevant items to be represented more strongly in posterior brain areas (notably visual
8 cortex), while prospectively relevant items were represented more strongly in frontal regions
9 (notably the Frontal Eye Fields). Under this scenario the representational overlap between current
10 and prospective items within the brain regions involved is expected to be minimal. Although crucial
11 for current theories of working memory, so far, studies have not directly compared the
12 representational pattern of current and prospective memories.

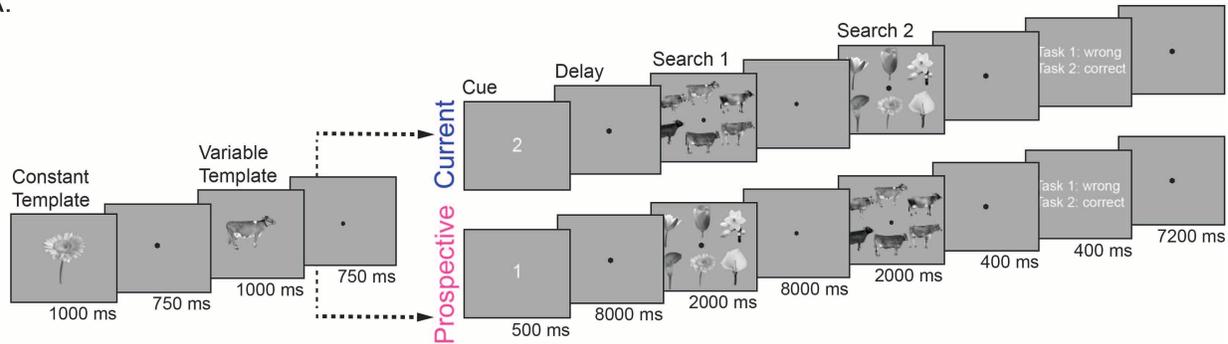
13 In two experiments we aimed to further understand how working memory distinguishes
14 between information relevant for either imminent or future goals. We asked observers to perform
15 two consecutive visual searches for particular target objects drawn from different object categories
16 (see Figure 1A and 1B for Experiments 1 and 2 respectively). Prior to search, these objects would
17 be maintained in working memory as target templates. Importantly, a cue indicated whether the
18 target template of interest would be relevant for the *first* search (turning it into a *current* template),
19 for the second search (turning it into a prospective template) or would not be relevant for either
20 search task (irrelevant condition, only in Experiment 2). Using MVPA of fMRI activity in object-
21 selective visual cortex, we directly compared the neural representations of these templates when
22 needed for the current search task, to when needed for the prospective search task. Experiment 1
23 served to establish the relationship between currently and prospectively relevant representations,
24 while Experiment 2 extended the comparison to representations that could be dropped from
25 memory entirely, as they became irrelevant for the subsequent tasks.

26 These experiments reveal a dissociation between currently relevant, prospectively relevant
27 and irrelevant templates based on category selective patterns in object-selective cortex. We find
28 that while observers are searching displays for the current target, the prospective search template
29 can nevertheless be temporarily decoded, extending the demonstration that prospective memories
30 can be reconstructed by sending unrelated activity through the network (Rose et al., 2016; Wolff et
31 al., 2017). Most importantly, we find that during the first search, the pattern of activity
32 corresponding to the prospective template is the inverse of the same template when it is currently
33 relevant. Thus, patterns reflecting prospective and current memory templates are systematically

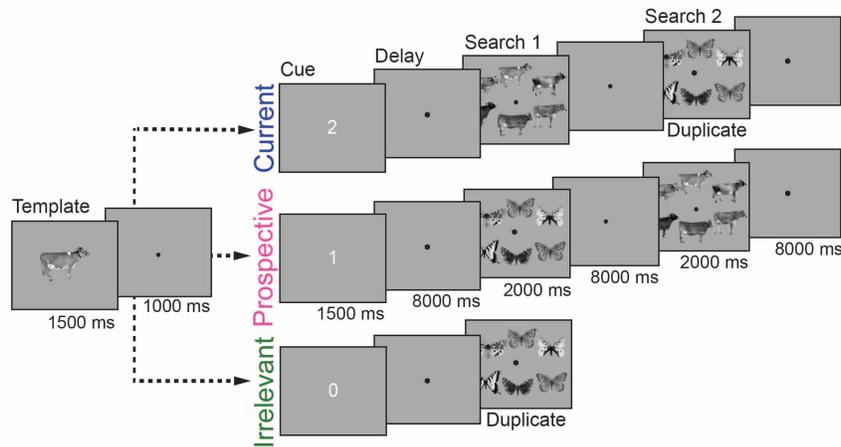
1 dissimilar, even when they belong to the same category. Experiment 2 further demonstrates that
 2 this inverse representational code is specific to the maintenance of information for future goals, as
 3 irrelevant representations did not show such an inversion. These results suggest that prospective
 4 templates are protected from interfering tasks by maintaining them in an opposite representational
 5 space.

6

A.



B.



7

8 **Figure 1. (A) Trial design for Experiment 1.** On each trial, participants performed two consecutive
 9 visual search tasks. The target objects for both search tasks were presented at the start of the trial.
 10 One of the objects could either be a cow, dresser or skate (variable template search; 4 exemplars per
 11 category), and was used for the decoding analyses. The other target was always the same flower
 12 (constant template search). The order of presentation (constant or variable template) was
 13 manipulated between trials, to create the two main conditions – one in which the variable template
 14 was currently relevant, the other in which it was prospectively relevant. To this end, a retro-cue (“1” or
 15 “2”) indicated which of the two previously memorized objects was the target in Search 1. The cue was
 16 followed by a delay, then the first search display, followed by a second delay and finally the second
 17 search display. Thus, in the Current condition, observers first searched for the variable template (cow,
 18 dresser, or skate), and then for the constant template (flower), while this order reversed in the
 19 Prospective condition. For each search display, participants indicated whether the target object was
 20 present or absent using a button press. At the end of each trial and run participants received feedback
 21 about their performance. **(B) Trial design for Experiment 2.** Here participants were presented with
 22 only one object (cow, dresser or skate) as the possible target template for one of two consecutive visual

1 *search tasks. Then a retro-cue appeared, when the cue was “1” the memorized object was a current*
2 *template, for Search 1; cue “2” indicated that the object was a prospective template, for Search 2;*
3 *finally, when the cue was “0” the memorized item was not a target in either search and thus it was*
4 *irrelevant in the trial. The remaining search task in Experiment 2 (either Search 2 in the Current*
5 *condition, or Search 1 in the Prospective and Irrelevant conditions) was a so-called duplicate search*
6 *task. In this task, butterflies, motorcycles or trees were presented and participants indicated whether*
7 *or not any one of the exemplars was shown twice in the display. Thus, here no search template could or*
8 *needed to be prepared. In the irrelevant condition participants only performed the duplicate search*
9 *task.*

10 **Source data 1.** Behavioral data for each participant of Experiment 1 and Experiment 2.

11 **Table supplement 1.** Behavioral data of Experiment 1.

12 **Table supplement 2.** Behavioral data of Experiment 2.

13 **Results**

14 **Experiment 1**

15 To examine the relationship between currently and prospectively relevant representations, on each
16 trial observers (N=24) performed two consecutive visual search tasks (Search 1 and Search 2). The
17 two to-be-sought-for target templates were presented at the start of each trial, after which a cue
18 indicated which of the two targets would have to be looked for first – thus making it currently
19 relevant, while the other target became prospectively relevant. To limit the working memory load,
20 and to maximize the chances of decoding current and prospective targets and their differences (See
21 Method section), we only varied the target from trial to trial for one of the two searches (thus
22 referred to as the “variable template search”). These targets served as the basis of the multivariate
23 pattern classification analyses, and could thus either be Current or Prospective in nature. The other
24 search was always for the same flower (referred to as the “constant template search”). The flower
25 search served as an additional task to assign current or prospective status to the variable template
26 (see methods), but the flower itself played no role in the classification analyses. For each search,
27 participants indicated whether the target object was present or absent among six exemplars of the
28 same category. Behavioral results show that both accuracy and search speed were better for the
29 constant template search than for the variable template search (see Behavioral Results Experiment
30 1 in Figure 1-table Supplement 1). Furthermore, the variable template search was more accurate
31 when performed first than when performed second, while order did not matter for the constant
32 template search performance. These results indicate that, as intended, working memory was indeed
33 involved more in the variable template than in the constant template. Our subsequent decoding
34 analyses are based on the variable template.

35

1 **fMRI results: Target template decoding as a function of current and prospective relevance.**

2 Our analyses targeted posterior fusiform cortex (pFs) which is known to be involved in
3 representing object categories, and which we independently mapped for each participant
4 (following Harel, Kravitz, & Baker, 2014; Kravitz, Saleem, Baker, Ungerleider, & Mishkin, 2013; Lee,
5 Kravitz, & Baker, 2013; Malach et al., 1995). To investigate whether we could decode memory
6 content for currently and prospectively relevant objects, we trained a classifier on the multivoxel
7 response patterns in pFs using each variable template category (i.e., cow, dresser and skate) for
8 each TR. Here we focus on the multivariate representation of the template categories of interest, as
9 shown in Figure 2, but the mean BOLD response for this area is also shown in Figure 1 - figure
10 supplement 1. First, we trained and tested the classifier separately for trials in which the target
11 category was currently relevant (during Search 1) and when the target category was prospectively
12 relevant (during Search 2, see Methods section for details). Object category classification
13 performance for this within-relevance decoding scheme is shown in Figure 2A. We focused our
14 statistical analysis on the averaged classification performance for three intervals in the trial (of
15 three TRs each; as predetermined on the basis of Lee et al. (2013), referred to as Delay, Search 1,
16 and Search 2; see Methods). We used paired t-tests ($N = 24$) to compare the classification
17 performance to chance (33.33%) for these intervals, as well as between Current and Prospective
18 conditions. Figure 2B shows the average activity for these time windows.

19 Second, while the within-relevance decoding scheme provides evidence for the presence of
20 current and prospective representations, it does not reveal whether these representations are
21 similar or different. Therefore, we additionally planned a cross-relevance decoding scheme in
22 which we trained the classifier when the objects were currently relevant and tested when the same
23 objects were prospectively relevant (referred to as PC), and vice versa (referred to as CP, see
24 Methods). Figures 2C and 2D show the classification accuracy for this cross-relevance decoding
25 scheme. Crucially, if current and prospective template representations are similar, above-chance
26 classification accuracy is expected. If representations are dissimilar in an unrelated fashion,
27 classification is expected to be at chance levels, while below-chance classification is predicted when
28 representations are dissimilar, but in a systematic, anti-correlated fashion. Our general starting
29 hypothesis was that while current and prospective representations would be similar during
30 encoding, they would become increasingly dissimilar during the course of the trial, due to reduced
31 activity or re-coding of the prospective item within the same network, while becoming similar again
32 when the prospective memories are revived for the second task.

33

1 ***The delay prior to the first search: Stronger decoding for current than for prospective***
2 ***templates, but similar representations.***

3 As can be seen in Figures 2A and 2B, during the Delay prior to search the within-relevance decoding
4 resulted in significant above chance object category decoding both when the variable template was
5 currently relevant ($t_{(1,23)} = 8.18, p < 0.001, d = 1.67$) and when prospectively relevant ($t_{(1,23)} = 5.67, p$
6 $< 0.001, d = 1.16$). However, while decoding accuracy for the current representation remained
7 significantly above chance up and beyond the Search 1 display, the prospective representation
8 returned to baseline during the delay. Notably, decoding performance was higher when the item
9 was currently relevant than when it was prospectively relevant (Current vs. Prospective: $t_{(1,23)} =$
10 $3.22, p = 0.004, d = 0.66$), consistent with its importance for the upcoming search task. Thus, object-
11 selective cortex proves sensitive to object category as well as task-relevance prior to search.

12 Next, we used the *cross-relevance* decoding scheme to assess whether current and
13 prospective targets shared the same neural representational pattern (see Figure 2B and 2C). This
14 analysis revealed strong above-chance classification of the template category, regardless of the
15 specific training scheme (PC: $t_{(1,23)} = 8.81, p < 0.001, d = 1.80$ or CP: $t_{(1,23)} = 9.04, p < 0.001, d = 1.85$).
16 There was no difference in decoding performance between the two schemes (PC vs. CP: $t_{(1,23)} = 1.43,$
17 $p = 0.167, d = 0.29$). These results indicate that during the delay prior to search, the
18 representational pattern of the category was similar regardless of the current or prospective status
19 of the object.

20 ***Search 1: The prospective template can be decoded during the first search, but is different from***
21 ***its current counterpart.*** Next, we wanted to know whether it was possible to successfully decode
22 the category of the prospective template while participants were searching for a different object.
23 As can be seen in in Figures 2A and 2B, in the Search 1 interval we observed clear decoding of the
24 object category when currently relevant (vs. 33.33%: $t_{(1,23)} = 11.57, p < 0.001, d = 2.36$), and this
25 was stronger than when the object was prospectively relevant (between-condition comparison:
26 $t_{(1,23)} = 9.42, p < 0.001, d = 1.92$). This is to be expected as during the first search of the Current
27 condition (i.e., variable template search) the current template category is actually presented on the
28 screen, whereas in the Prospective condition the objects on screen (i.e., flowers) differ from the
29 prospective target in memory. Importantly, we were still able to also decode the prospective
30 category during the first search (vs. 33.33%; $t_{(1,23)} = 1.90, p = 0.035, d = 0.39$). In other words, the
31 prospective category re-emerges when observers are actively searching for an unrelated target.

32 This then raises the question as to whether the re-emerging prospective representation
33 resembles its counterpart when currently relevant. To assess this we used the cross-relevance

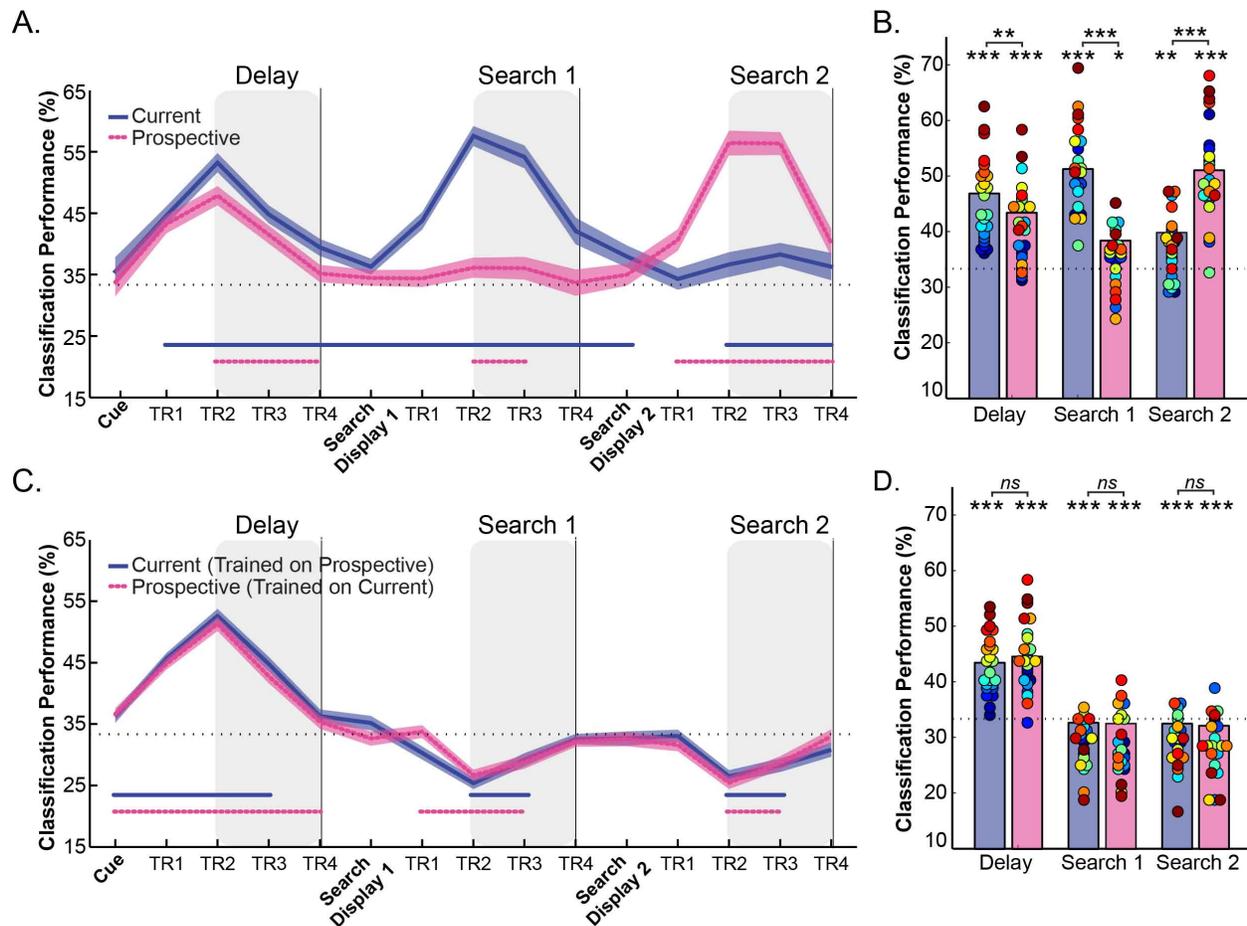
1 decoding scheme. Remarkably, here we observed *below*-chance decoding performance during
2 Search 1 (CP: $t_{(1,23)} = -4.79$, $p < 0.001$, $d = -0.98$, and PC: $t_{(1,23)} = -3.67$, $p = 0.001$, $d = -0.75$). Although
3 we hypothesized that current and prospective templates might differ in representational format,
4 and thus cross-relevance decoding accuracy might have been reduced, we did not expect the sign of
5 decoding to flip. The reliable deviation from chance further confirms that information on the
6 prospective memory was present in object-selective cortex during the first search. In addition, the
7 fact that decoding was below chance suggests that current and prospective representations of the
8 same object category were represented through opposite multivariate patterns.

9 Finally, we assessed how the dissociation between current and prospective representations
10 generalized across the different phases of the trial. As Figure 1-figure supplement 2 shows, the
11 pattern of activity prior to search is very similar to that during search for currently relevant
12 representations, whereas prospectively relevant representations during the first search are
13 markedly dissimilar from the same categories during the delay period prior to search. They then
14 become similar again when retrieved for Search 2. Thus, while currently relevant representations
15 remained constant from delay to search, the prospective representation was transformed from
16 being similarly represented prior to search to being differently represented during search for the
17 currently relevant item, followed by a reactivation for the second task.

18 ***Search 2: Decoding of the first, now no longer relevant target during the second search.***

19 Although not the primary goal of our study, we conducted the same analyses also for Search 2. As
20 expected, here we saw the pattern reverse (see Figure 2A and 2B). In the within-relevance decoding
21 scheme, we observed strong decoding of the category of the prospective target, which by now had
22 become currently task-relevant (against chance, 33%: $t_{(1,23)} = 9.86$, $p < 0.001$, $d = 2.01$). This
23 decoding was stronger than for the previously current search target, which was now no longer
24 relevant ($t_{(1,23)} = -7.51$, $p < 0.001$, $d = -1.53$). Nevertheless, and unexpectedly, we also observed
25 above-chance decoding for this first target during Search 2 ($t_{(1,23)} = 3.30$, $p = 0.002$, $d = 0.67$, blue).
26 Note that this reflects classification of a target that is no *longer* relevant, whereas during Search 1 it
27 reflected the target that was not *yet* relevant. Moreover, the cross-relevance decoding scheme also
28 shows a pattern similar to what was observed during Search 1 (Figure 2C and 2D; and see also
29 Figure S2 for the generalization across time). We found below-chance decoding in the same time
30 range for both classification schemes (CP: $t_{(1,23)} = -4.65$, $p < 0.001$, $d = -0.95$ and PC: $t_{(1,23)} = -4.49$, $p <$
31 0.001 , $d = -0.92$). We will return to the re-emergence of the no longer relevant target later.

33



1
 2 **Figure 2. Within-relevance and Cross-relevance object category decoding in pFs.** (A) Time
 3 **course of the Within-relevance decoding** where the classifier was trained and tested either within
 4 the current, or within the prospective conditions and (B) **Average decoding accuracy** within the time
 5 intervals shown by the shaded areas in (A). Decoding accuracy was higher for currently relevant
 6 templates (blue) than for prospectively relevant templates (pink) during the Delay and Search 1
 7 intervals, and vice versa during the Search 2 interval. At the same time, the prospective template could
 8 still be reliably decoded during the first search, while the no longer relevant target could be decoded
 9 during the second search. (C) **Time course of the Cross-relevance category decoding** where the
 10 classifier was trained on current relevance, tested on prospective relevance, or vice versa and (D)
 11 **Average decoding accuracy** within the time intervals as shaded in (C). This resulted in above-
 12 chance decoding during the Delay prior to search (suggesting similar representations for current and
 13 prospective templates) but below-chance decoding during Search 1 and Search 2 (suggesting partially
 14 opposite representations). Shaded blue and pink areas indicate within-subjects s.e.m. Blue and pink
 15 horizontal lines at the bottom of the line graphs indicate time points that significantly differ from
 16 chance ($p < 0.05$, uncorrected). In the bar-plots, colored dots indicate individual participant data, $N =$
 17 24. * $p < 0.05$, ** $p < .01$, *** $p < .001$, ns: not significant.

18 **Source data 1.** Decoding performance for each participant of Experiment 1. Includes source code to
 19 perform statistical analysis and produce Figure 2.

20 **Figure supplement 1.** Mean BOLD response of Experiment 1.

21 **Source data 2.** Mean BOLD response for each participant. Includes source code to perform statistical
 22 analysis and produce Figure 2 supplement 1.

1 **Figure supplement 2.** *Cross-temporal generalization matrices for object category decoding as a*
2 *function of Relevance*

3 **Source data 3.** *Cross-temporal generalization matrices for each participant Includes source code to*
4 *perform statistical analysis and produce Figure 2 supplement 2.*
5

6 ***Representational dissimilarity analyses comparing current and prospective representations.***

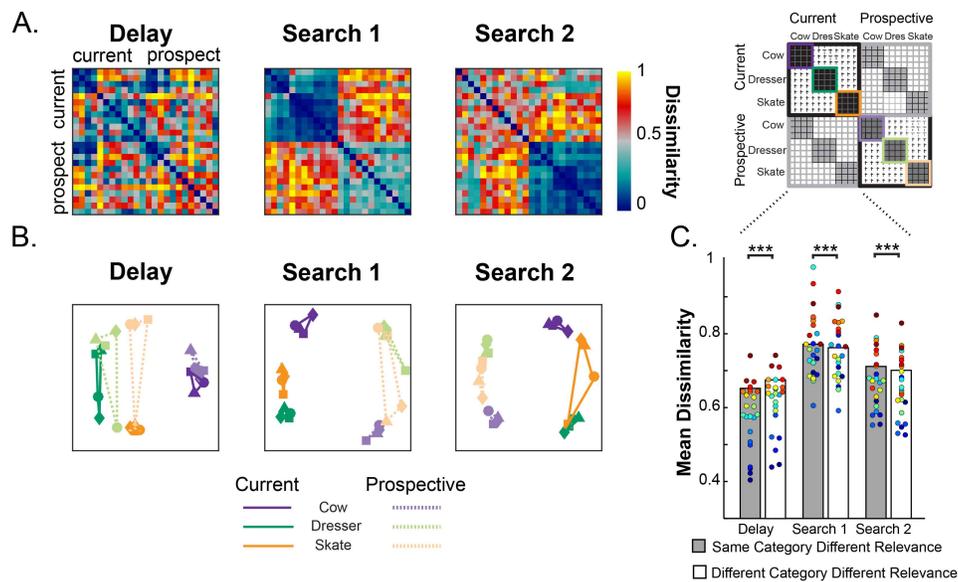
7 To further elucidate the relationship between current and prospective representations, Figure 3A
8 shows, for each interval of interest (Delay, Search 1, Search 2), the representational dissimilarity
9 matrices (Kriegeskorte & Kievit, 2013; Kriegeskorte et al., 2008). Specifically, we cross-correlated
10 the voxel response patterns for every possible pair of the 24 stimulus/relevance combinations (4
11 exemplars x 3 categories (Cow, Dresser and Skate) x 2 relevance (Current and Prospective; see
12 Figure 3A right panel and Methods for further details). To further visualize the relative position in
13 multivariate space of each object category as a function of relevance, Figure 3B shows
14 multidimensional scaling (MDS) graphs: the shorter the distance between categories the greater the
15 representational similarity across voxels.

16 As can be seen from Figures 3A and 3B, throughout the course of the trial the neural
17 representations moved from predominantly category space during the Delay period prior to search
18 to predominantly relevance space during the two searches. Prior to search, objects grouped largely
19 according to category, irrespective of relevance. This confirms that currently relevant and
20 prospectively relevant objects were initially represented in similar ways in pFs. During Search 1, a
21 clear relevance-driven distinction emerged between the neural object category representations.
22 Note that this overall effect of relevance is probably driven by the fact that during search the
23 currently relevant object category was presented in the display, while in the prospective condition
24 the unrelated (flower) displays were presented. More interesting though is the finding that
25 currently and prospectively relevant objects from the same category were represented as the most
26 *dissimilar*, as is illustrated by the representations taking opposite corners in the MDS plot in Figure
27 3B. For example, while all four exemplars of the cow category clustered together when all current,
28 or all prospective, current cows were most separated from prospective cows – to the extent that
29 current cows were more similar to prospective skates and dressers than they were to prospective
30 cows. The same pattern held for the two other categories.

31 To statistically test these effects, we computed the average dissimilarity between current
32 and prospective objects, separately for when drawn from the same category (e.g. current cow
33 versus prospective cow) and when drawn from a different category (e.g. current cow versus
34 prospective skate/dresser) and used paired t-tests (N=24). Figure 3C shows these average same
35 and different category dissimilarity values across relevance. During the Delay prior to the first

1 search, as expected, *same* category representations were more similar than *different* category
 2 representations across relevance ($t_{(1,23)} = -5.82, p < 0.001, d = -1.28$, as performed on Fisher-
 3 transformed 1-r values). In contrast, during Search 1, prospective targets differed most from
 4 current targets when they belonged to the same category, more so than when they belonged to
 5 different categories ($t_{(1,23)} = 3.06, p = 0.005, d = 0.64$). Likewise, during Search 2, no longer relevant
 6 targets were less similar from relevant targets when they belonged to the same category, than
 7 when they belonged to different categories ($t_{(1,23)} = 4.75, p < 0.001, d = 0.97$). Thus these analyses
 8 statistically confirm what we can observe from the MDS plots, namely that current and prospective
 9 objects move from similar to opposite representations.

10
 11



12

13 **Figure 3. (A) Representational dissimilarity of object representations in pFs.** Representational
 14 dissimilarity matrices for the different variable template categories during Delay, Search 1 and Search
 15 2, as a function of relevance (current and prospective). Blue indicates that representations are more
 16 similar while red indicates more dissimilar **(B) Multidimensional scaling plots** of the same similarity
 17 values, for the same Delay, Search 1 and Search 2 intervals. The four exemplars within each category
 18 are represented with different shapes (squares, triangles, circles and diamonds). The closer in space
 19 the more similar the neural representations. As a trial unfolds, object representations move from
 20 predominantly object category space during the delay prior to search (e.g. a cow) into predominantly
 21 relevance space (e.g. current target) during search, where current and prospective targets of the same
 22 category are represented by partly opposite representational patterns. **(C) Comparing dissimilarity**
 23 **between Current and Prospective** items when they are drawn from the same category versus when
 24 they are drawn from different categories. Representations prior to search within the same category
 25 are more similar than different categories, but this reverses during the searches. Colored dots indicate
 26 individual participant data, ** $p < .01$, *** $p < .001$.

27 **Source data 1.** RDM for each participant of Experiment 1. Includes source code to perform statistical
 28 analysis and produce Figure 3.

1 **Experiment 2**

2 What causes current and prospective representations to anti-correlate? One possibility is that the
3 brain separates current from prospective templates within the same neuronal ensembles by
4 actively transforming the representational pattern of prospective templates to be opposite to that
5 of current templates. A second possibility is that the mechanism of making a representation
6 prospective is more passive, and that the reversed pattern results from simply temporarily
7 dropping a representation from memory, as it is temporarily irrelevant. One piece of evidence from
8 Experiment 1 suggests that postponing and dropping an item may indeed involve similar
9 mechanisms: During Search 2 we observed similar negative decoding for target representations
10 that had served the first search and that were thus no longer necessary. However, getting rid of a
11 competing target representation when switching to a new search may also involve active
12 mechanisms, now to prevent interference from the past target rather than from a future target.
13 Results from our lab indeed suggest such a suppression of previous targets (de Vries, van Driel,
14 Karacaoglu, & Olivers, 2018). A third possibility is that the pattern of reversal has little to do with
15 memory whatsoever, but is simply a remnant of stimulus-related activity during encoding.
16 Specifically, presenting the to-be-memorized stimulus may result in neural adaptation (Henson &
17 Rugg, 2003; Larsson & Smith, 2012; Vautin & Berkley, 1977) or in a BOLD-related undershoot
18 (Huettel & McCarthy, 2000), each of which would predict a reduced voxel response to later
19 stimulation. Finally, the observed pattern of Experiment 1 may have been caused by certain
20 idiosyncrasies of the experiment, most notably the fact that we always used the same flower target
21 in the constant template condition, which may have led to either stimulus-specific or overall task
22 difficulty-related interactions.

23 To address this, Experiment 2 sought to replicate and extend the main findings with a
24 number of design changes. The most important change was the inclusion of a third condition, in
25 which the memory item was cued to be irrelevant for any of the search tasks (see Figure 1B and
26 Methods for details). Importantly, this Irrelevant condition matched the Prospective condition in all
27 aspects – including the visual input – up to and including the first search, making the initial sensory
28 processing identical across conditions. However, while in the Prospective condition the memory
29 item was cued to become relevant only after the first search, in the Irrelevant condition the memory
30 item was cued to become irrelevant altogether. Therefore, any differences in results across the
31 Irrelevant and Prospective conditions can only be attributed to the future relevance of the
32 memorized item and not to any systematic differences between memory items or search displays.

1 For the same reason, neither can any differences between irrelevant and prospective
2 representations be attributed to passive, sensory-related adaptation or BOLD undershoot.

3 Furthermore, we simplified the design by keeping the variable template search (here
4 referred to as simply “template search”), but replacing the constant template task of Experiment 1
5 with what we call a “duplicate search” task, in which participants indicated whether or not one of
6 the objects in the search display appeared twice (see Figure 1B and Methods for details, as well as
7 de Vries et al., 2018). Note that such duplicate search does not require a template, because all the
8 information needed to perform the task is in the search display itself. At the same time, it still
9 engenders a dissociation between current and prospective memory templates over time. As a
10 result, observers only had to remember a single target template per trial, which was either the
11 target for the first search (Current; with the second search being a duplicate search), or the target
12 for the second search (Prospective; with the first search being the duplicate search), or was deemed
13 irrelevant after all (Irrelevant condition; with the first and only search being a duplicate search).
14 Finally, just like the stimuli for the template search were drawn from three different categories
15 (cows, skates, and dressers), we varied the stimuli in the duplicate search such that they were also
16 drawn from three different categories (specifically butterflies, motorcycles, and trees), in order to
17 assess whether the (below-chance) decoding of prospective representations during search
18 generalizes across a range of different categories. The behavioral results show that performance in
19 the template search and the duplicate search were comparable in terms of accuracy, and that the
20 template search was actually faster than the duplicate search (see and Figure 1-table Supplement
21 2). Thus, here overall task difficulty if anything showed the reverse pattern compared to
22 Experiment 1. Yet, the fMRI findings show a pattern very similar to that of Experiment 1, as
23 reported next.

24 **fMRI results: Target category decoding as a function of task relevance.**

25 As in Experiment 1, here we focus on multivariate representational patterns within the same three
26 intervals: Delay, Search 1, and Search 2 (see Figure 4 and Methods for details), here for Current,
27 Prospective, and Irrelevant objects. The mean overall BOLD response in pFS is shown in Figure 4-
28 figure supplement 1.

29

30 ***The delay prior to the first search: Stronger decoding for current and prospective templates*** 31 ***than for irrelevant items, but similar representations across conditions.***

32 Figures 4A and 4B show the decoding accuracy during the Delay prior to search. A one-way ANOVA

1 on the within-relevance decoding of Current, Prospective and Irrelevant conditions revealed a
2 significant effect of condition ($F_{(2,48)} = 4.40, p = 0.018, \eta_p^2 = 0.15$). As shown in Figure 4B, decoding
3 accuracy was highest for Current, lowest for Irrelevant, with the Prospective condition in between.
4 There was significant above-chance object category decoding for all relevance conditions: Current
5 ($t_{(1,24)} = 10.76, p < 0.001, d = 2.15$) Prospective ($t_{(1,24)} = 5.88, p < 0.001, d = 1.17$) and Irrelevant
6 ($t_{(1,24)} = 5.61, p < 0.001, d = 1.12$). Pairwise comparisons revealed the difference between Current
7 and Irrelevant to be significant ($t_{(1,24)} = 3.74, p = 0.001, d = 0.74$). In contrast to Experiment 1
8 though, the difference in decoding accuracy between the Current and Prospective conditions was
9 not significant ($t_{(1,24)} = 1.33, p = 0.193, d = 0.26$), nor was there a significant difference between the
10 Prospective and Irrelevant conditions ($t_{(1,24)} = 1.39, p = 0.175, d = 0.27$). We will return to the
11 possible reasons for this difference between experiments in the General Discussion.

12 Next, we used the *cross-relevance* decoding scheme to assess whether current, prospective
13 and irrelevant items shared the same neural representational pattern (see Figure 4C and 4D). The
14 classification performance of each particular train-test scheme and its converse counterpart were
15 averaged (see Methods for details). This analysis revealed above-chance classification for the three
16 cross-relevance decoding schemes: Current-Prospective ($t_{(1,24)} = 9.85, p < 0.001, d = 1.97$), Current-
17 Irrelevant ($t_{(1,24)} = 8.10, p < 0.001, d = 1.62$) and Prospective-Irrelevant ($t_{(1,24)} = 6.90, p < 0.001, d =$
18 1.38). A one-way repeated measures ANOVA with decoding scheme as factor revealed no significant
19 differences between decoding schemes ($F_{(2,48)} = 2.4, p = 0.101, \eta_p^2 = 0.09$). Taken together, these
20 results indicate that during the delay prior to search, the representational pattern of the memorized
21 object category was similar regardless of the current, prospective or irrelevant status of the object.

22 23 ***Search 1: The prospective target but not the irrelevant item can be decoded during the first*** 24 ***search.***

25 A one-way ANOVA on the within-relevance decoding during Search 1 showed a significant
26 difference in decoding accuracy ($F_{(1.6,38.77)} = 35.17, p < 0.001, \eta_p^2 = 0.59$, corrected for sphericity).
27 Decoding accuracy for the current template was above chance ($t_{(1,24)} = 9.47, p < 0.001, d = 1.89$) and
28 significantly higher than for both the prospective template ($t_{(1,24)} = 5.53, p < 0.001, d = 1.10$) and
29 the irrelevant item ($t_{(1,24)} = 7.29, p < 0.001, d = 1.45$). This was to be expected since during the first
30 search of the Current condition, the current template category was on the screen. More importantly,
31 and in line with Experiment 1, we were able to decode the prospective category at above-chance
32 levels during Search 1 ($t_{(1,24)} = 3.35, p = 0.003, d = 0.67$). At the same time, decoding of the *Irrelevant*

1 category remained at chance ($t_{(1,24)} = 0.60, p = 0.550, d = 0.12$), and significantly weaker than for
2 Prospective condition ($t_{(1,24)} = 2.27, p = 0.032, d = 0.45$). Thus, information about the prospective
3 template can be recovered while participants perform a different search task, whereas completely
4 irrelevant categories are no longer decodable.

5 Next, we used the cross-relevance decoding scheme to evaluate whether the
6 representations of the prospective targets and irrelevant items resemble their current counterpart
7 (see Figure 4C and 4D). A one-way ANOVA with decoding scheme (Current-Prospective, Current-
8 Irrelevant, and Prospective-Irrelevant) revealed a reliable effect ($F_{(2,48)} = 16.35, p < 0.001, \eta_p^2 =$
9 0.40). Replicating Experiment 1, we observed strong *below*-chance decoding during Search 1 for the
10 Current-Prospective cross-relevance classification ($t_{(1,24)} = -5.92, p < 0.001, d = -1.18$), while the
11 Current-Irrelevant ($t_{(1,24)} = -1.89, p = 0.071, d = -0.37$) and the Prospective-Irrelevant ($t_{(1,24)} = 1.99,$
12 $p = 0.058, d = 0.39$) cross-decoding schemes did not significantly differ from chance. Most
13 importantly, decoding accuracy for the Current-Prospective scheme was significantly *lower* than for
14 the Current-Irrelevant scheme ($t_{(1,24)} = -2.95, p = 0.007, d = -0.59$), indicating that the Prospective
15 condition involved a stronger representational transformation than the Irrelevant condition. Since
16 the prospective and irrelevant trials contained exactly the same visual input (see Methods), this
17 result suggests that this transformation is driven by the future relevance of the prospective
18 template rather than by some passive, automatic mechanism that would be the same for
19 prospective and irrelevant representations (such as BOLD undershoot, see General Discussion).

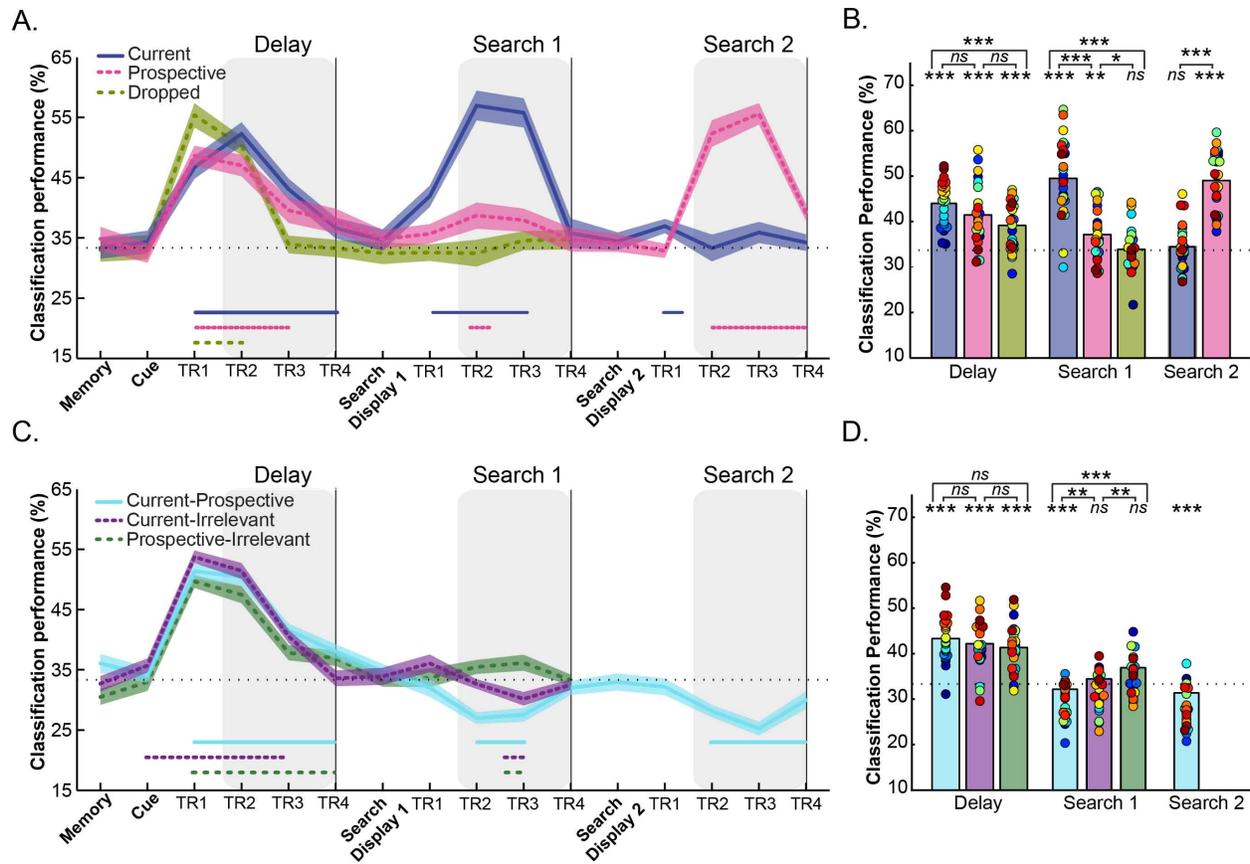
20

21 ***Search 2: Evidence for the first, no-longer relevant target during the second search.***

22 In line with Experiment 1, the pattern reversed for Search 2 (see Figure 4). In the within-relevance
23 decoding scheme, we observed strong decoding of the category of the prospectively relevant target,
24 which by now had become task-relevant ($t_{(1,24)} = 12.54, p < 0.001, d = 2.50$), and decoding of the
25 Prospective -now relevant- category was stronger than for the previously *current* -now no longer
26 relevant- search target ($t_{(1,23)} = -8.84, p < 0.001, d = -1.76$). In contrast to Experiment 1, where we
27 unexpectedly observed above-chance decoding for this no longer relevant target during Search 2,
28 here the current condition was at chance ($t_{(1,24)} = 1.14, p = 0.265, d = 0.22$). However, similar to
29 Experiment 1, we found below-chance decoding in the same time range for the Current-Prospective
30 cross-relevance decoding scheme (vs. 33.3%: $t_{(1,24)} = -6.57, p < 0.001, d = -1.31$), indicating that there
31 was information present on the previously relevant target.

32

33



1
 2 **Figure 4. Within-relevance and cross-relevance object category decoding in pFs. (A) Time**
 3 **course of within-relevance decoding and (B) Averaged decoding accuracy** within the time
 4 intervals shown by the shaded areas in A. During the delay, object category decoding was higher for
 5 currently relevant objects (blue) than for irrelevant objects (green) with in between decoding
 6 accuracy for prospective templates. During Search 1 the current template showed higher decoding
 7 accuracy than the prospective template and the irrelevant item. Importantly, the category of the
 8 prospective template could also be decoded during the first search, while the irrelevant category was
 9 at chance. During Search 2 the prospective (now current) category was clearly decodable while the
 10 formerly current (now no longer relevant) category was at chance **(C) Time course of cross-**
 11 **relevance decoding and (D) Averaged decoding accuracy** within the time intervals shown by the
 12 shaded areas in C. Classification was above chance for all decoding combinations during the Delay
 13 prior to search, suggesting similar representations for current, prospective and irrelevant objects. In
 14 contrast, we observed below-chance decoding during Search 1 and Search 2 for the Current-
 15 Prospective (blue) cross-relevance scheme (suggesting systematically different representations);
 16 importantly, this was stronger than for Current-Irrelevant (purple; during Search 1). Current-
 17 Irrelevant (purple) and Prospective-Irrelevant (green) cross-classification schemes were at chance.
 18 Shaded areas indicate within -subjects s.e.m. Blue and pink, purple and green horizontal lines at the
 19 bottom of the line graphs indicate time points that significantly differ from chance ($p < 0.05$,
 20 uncorrected). In the bar-plots, colored dots indicate individual participant data, $N = 25$. * $p < 0.05$, ** $p <$
 21 $.01$, *** $p < .001$, ns: not significant.

22 **Source data 1.** Decoding performance for each participant of Experiment 2. Includes source code to
 23 perform statistical analysis and produce Figure 4.

24 **Figure supplement 1.** Mean BOLD response of Experiment 2.

1 **Source data 2.** Mean BOLD response for each participant of Experiment 2. Includes source code to
2 perform statistical analysis and produce Figure 4 supplement 1.

3
4 **Representational dissimilarity analyses comparing current, prospective and irrelevant**
5 **representations.**

6 Next, we created representational dissimilarity matrices and multidimensional scaling (MDS)
7 graphs for each interval of interest (Delay, Search 1, Search 2) for the three possible condition
8 combinations: Current-Prospective, Current-Irrelevant and Prospective-Irrelevant (see Figure 5).
9 First, as can be seen in Figure 5A, 5B and 5C, we replicated the results from Experiment 1. In the
10 Delay period, the neural representations of current and prospective targets of the *same* category
11 were more similar (i.e., less dissimilarity) than the representations of targets from *different*
12 categories ($t_{(1,24)} = -9.14, p < 0.001, d = -1.82$, Figure 5C). In contrast, during Search 1 and Search 2,
13 currently and prospectively relevant objects from the same category were more *dissimilar*: As in
14 Experiment 1, prospective targets differed most from current targets when they belonged to the
15 *same* category than when they belonged to *different* categories (Search 1: $t_{(1,24)} = 5.20, p < 0.001, d =$
16 1.04 and Search 2: $t_{(1,24)} = 6.45, p < 0.001, d = 1.29$). Again, throughout the course of the trial, the
17 neural representations of the target objects moved from predominantly category space in the Delay
18 period prior to search to predominantly relevance space during the two searches, where current
19 and prospective objects were represented in opposite corners of the representational space.

20 The results for the Current-Irrelevant (Figures 5D, 5E and 5F) and Prospective-Irrelevant
21 (Figures 5G, 5H and 5I) condition combinations were very comparable during the delay prior to the
22 first search. Here too, the neural representation of targets of the same category were more similar
23 than when they were from different categories (Current-Irrelevant: $t_{(1,24)} = -8.27, p < 0.001, d = -$
24 1.65 and Prospective-Irrelevant: $t_{(1,24)} = -4.73, p = 0.001, d = -0.94$). Importantly however, and in
25 contrast to Prospective items, for Irrelevant items we did not find evidence that the representations
26 were warped in opposite corners of the multivariate space during Search 1, as the (dis)similarity
27 between targets belonging to the same or to different categories was equal for the Current-
28 Irrelevant ($t_{(1,24)} = 0.94, p = 0.354, d = 0.89$) and Prospective-Irrelevant ($t_{(1,24)} = -1.36, p = 0.184, d$
29 $= -0.27$) comparisons.

30 The uniqueness of the prospective transformation was further confirmed by directly
31 comparing the relative similarity difference (by subtracting the mean target dissimilarity when of
32 the *same* category from the mean dissimilarity when of *different* category) for each relevance
33 combination (i.e., Current-Prospective, Current-Irrelevant and Prospective-Irrelevant). This
34 showed that same category representations were indeed relatively more dissimilar for Current-

1 Prospective comparisons than for Current-Irrelevant ($t_{(1,24)} = 2.88, p = 0.008, d = 0.57$) and
2 Prospective-Irrelevant ($t_{(1,24)} = 4.69, p < 0.001, d = 0.93$) comparisons.

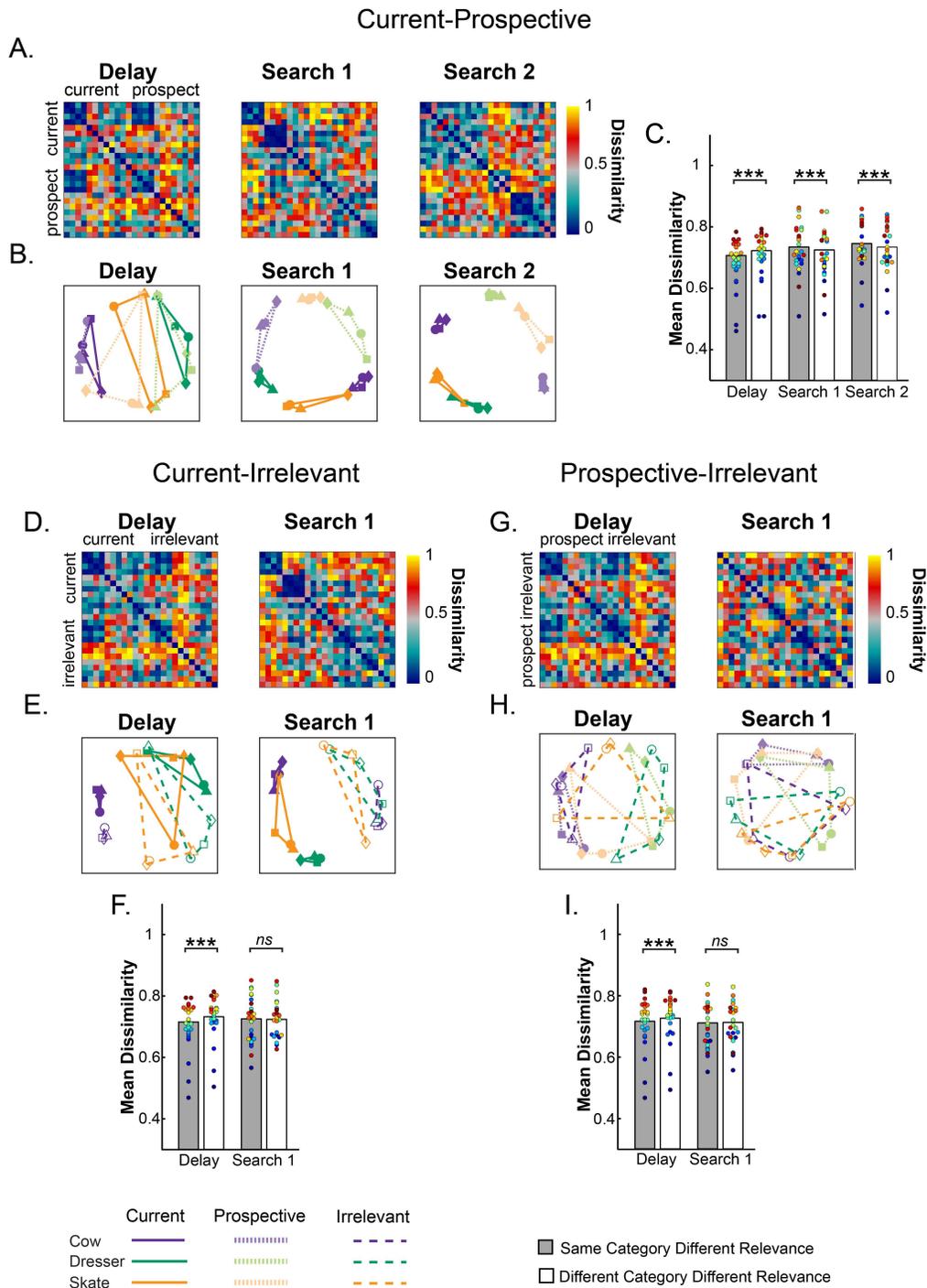
3 Taken together, the results of Experiment 2 confirm the most important findings of
4 Experiment 1. We find a pronounced anti-correlated representation for prospective targets
5 compared to the same targets when current. Experiment 2 furthermore shows that it is not the
6 temporary irrelevance, but the prospective relevance that causes this transformation, because little
7 to no anti-correlation was found for items that were irrelevant altogether. Since prospective and
8 irrelevant trials contained exactly the same visual input (up to Search 1 inclusive), the
9 transformation of prospective representations cannot be the result of basic stimulus-induced
10 neural adaptation or BOLD response properties.

11

12

13

14



1

2 **Figure 5. Representational dissimilarity analysis of object representations in pFs (A,D,G)**
 3 **Representational dissimilarity matrices and (B,E,H) Multidimensional scaling plots of the same**
 4 **similarity values for the different target object categories during Delay, Search 1 and Search 2, as a**
 5 **function of relevance. The four exemplars within each category are represented with different shapes**
 6 **(squares, triangles, circles and diamonds). (A,B) Comparing Current to Prospective: Object**
 7 **representations moved from predominantly object category space (e.g. a cow) during the Delay period**
 8 **to predominantly relevance space during search, where current and prospective targets of the same**
 9 **category were represented by partly opposite representational patterns. (C) Current-Prospective**

1 **Mean Dissimilarity.** During the delay prior to search, representations within the same category are
2 more similar than of different categories; however, the pattern reverses during Search 1 and Search 2,
3 where representations within the same category are more dissimilar than of different categories. **(D,**
4 **E, F) Comparing Current to Irrelevant and (G, H, I) Comparing Prospective to Irrelevant.** During
5 the delay prior to search, targets of the same categories were more similar than of different
6 categories, with no opposite representational pattern during search. Colored dots indicate individual
7 participant data, ** $p < .01$, *** $p < .001$, ns: not significant.
8 **Source data 1.** RDM for each participant of Experiment 2. Includes source code to perform statistical
9 analysis and produce Figure 5.

10

11

General Discussion

12 It has been shown that the content of working memory can be decoded from multivariate patterns
13 of voxel activity when observers are remembering an item for a single task (Albers, Kok, Toni,
14 Dijkerman, & de Lange, 2013; Harrison & Tong, 2009; Lewis-Peacock & Postle, 2012; Serences,
15 Ester, Vogel, & Awh, 2009). Furthermore, working memory representations have been shown to
16 adapt to the specific task goal, as the representation of the same object changes depending on the
17 nature of the upcoming test (Lee et al., 2013; Myers, Stokes, & Nobre, 2017). Here we provide
18 further evidence for the flexibility of working memory by showing how it distinguishes between
19 relevant and irrelevant representations on the one hand, and between relevant representations for
20 current and future task goals on the other hand, as representations adapt to the order in which they
21 are required in multiple task sequences.

22 In line with earlier work in different task and stimulus domains (LaRocque et al., 2013;
23 2017; Lewis-Peacock & Postle, 2012; Wolff et al., 2017), we observed that objects required for an
24 upcoming search task are represented more strongly than when the same objects are only used
25 prospectively (Experiment 1) or when they are completely irrelevant (Experiment 2), as was
26 indicated by stronger classification performance in object-selective cortex prior to the first search.
27 Furthermore, corroborating findings by Rose and colleagues (2016), we observed that the
28 prospectively relevant memory could be reconstructed during task-irrelevant stimulation, here
29 during the first search for an unrelated stimulus. This was shown in two ways: First, both
30 experiments showed above-chance classification of the prospective item during the first search,
31 whereas Experiment 2 showed that classifier performance for irrelevant representations remained
32 at chance and was significantly worse than for prospective representations. Second, using a cross-
33 relevance training scheme, where the classifier was trained when the item was current and tested
34 when it was prospective (or vice versa), we also found in both experiments that decoding
35 performance reliably differed from chance – but now in a negative direction.

1 Our results are the first to reveal a direct relationship between currently relevant object
2 representations on the one hand, and prospective representations on the other. Prior to the first
3 search current, prospective and irrelevant representations were very similar, as cross-relevance
4 decoding (training on one status while testing on the other) showed above-chance performance for
5 all classification schemes and there were no differences across conditions. However, while the
6 representation for prospective objects clearly reversed during search, such a reversal was weak to
7 absent for irrelevant objects. We observed that during search, the reconstructed prospective
8 memory representations of objects proved *dissimilar* from their current counterparts. Importantly,
9 they differed in a systematic manner, to the extent that prospective representations were even
10 more dissimilar from current representations of the *same* object category than representations of a
11 *different* object category, and were characterized by an inverse correlation with current
12 representations. Conversely, the initial similarity between current and irrelevant targets, which
13 was evident during the delay, mostly disappeared during search leading to the same dissimilarity
14 values for targets of the same vs. different category. Thus, while irrelevant targets simply appeared
15 to decay from memory, prospective target memories were transformed. We point out that similar
16 (as yet unpublished) findings have recently been reported by (Yu & Postle, n.d.), for a different
17 stimulus class and a different brain area. They asked participants to memorize two oriented
18 gratings, one for a first test (making it current), the other for a second test (making it prospective).
19 Using multivariate inverted encoding models (IEM) of occipital cortex activity, they could
20 reconstruct prospectively relevant orientations with models trained on currently relevant
21 orientations, but here too a reversed pattern emerged. Taken together, both our and these findings
22 indicate that prospective targets may be dissociated from current targets in two ways. First, they
23 appear distinct in that current representations are activity-based, whereas prospective
24 representations are responsivity-based. Second, prospective targets are represented through a
25 responsivity pattern different to that of current target activity, where the most active part becomes
26 the relatively least responsive and vice versa.

27 An issue left unresolved in the present study is whether the observed transformation of
28 prospective representations extends to specific exemplars. That is, are current and prospective
29 representations even more dissimilar when representing the same exemplar? Unfortunately, the
30 limited number of trials per exemplar and condition precluded a useful analysis here, and future
31 studies should directly test the item-specificity of the transformation of prospective
32 representations.

1 Imaging studies of visual attention have demonstrated that assigning attentional priority to
2 task-relevant objects at the expense of irrelevant objects leads to the transformation of the
3 representational space, by relatively enhancing target-related distributed activity (Reddy,
4 Kanwisher, & VanRullen, 2009) and by recruiting additional resources as neurons shift their tuning
5 towards the attended object category (Çukur, Nishimoto, Huth, & Gallant, 2013; Nastase et al.,
6 2017). Here, such relative enhancement can explain the pattern of results during the delay period
7 prior to the first search task in Experiment 1, where we found a difference in representational
8 strength between current and prospective templates, while their representations remained similar.
9 Note that no such difference occurred in the delay period of Experiment 2. This discrepancy
10 between experiments may be explained by assuming a direct competition between current and
11 prospective templates, which was only present in Experiment 1 (where there was a variable and a
12 constant template). In Experiment 2 there was only one object to remember, allowing observers to
13 devote the same resources in all conditions.

14 A crucial question that our data does not answer is what the exact mechanism is behind the
15 transformation from current to prospective representations. One possibility is the involvement of
16 an active cognitive control mechanism, which specifically attempts to dissociate prospectively
17 relevant from currently relevant representations in order to prevent task interference. Such control
18 mechanisms might be exerted through feedback connections emanating from frontal areas central
19 to counteracting unwanted or task-irrelevant information (Anderson et al., 2004; Banich,
20 Mackiewicz Seghete, Depue, & Burgess, 2015; de Vries, et al., 2018; Depue, Curran, & Banich, 2007;
21 Reeder, Olivers, & Pollmann, 2017). Interestingly, an earlier study of memory retrieval has shown
22 suppression of voxel patterns in ventral object-related cortex which were associated with task-
23 irrelevant memories of learned object pictures, leading to comparable patterns of representational
24 dissimilarity as here (Wimber, Alink, Charest, Kriegeskorte, & Anderson, 2015). Initial evidence for
25 the suppression of temporarily irrelevant items also comes from a study by Peters and colleagues
26 (2012), who used a similar task design as ours. They asked observers to consecutively look for a
27 particular house target and a particular face target (or vice versa) in rapid streams of house/face
28 distractors. They found the overall BOLD signal to be reduced in either house or face selective areas
29 in response to house/face stimuli when the respective target was prospectively relevant. Here we
30 show how changing task relevance within working memory specifically affects the cortical pattern
31 of activation within memory while observers perform a different search task. In this respect it is
32 also interesting to note that in both experiments we found evidence for an inversion both for
33 temporarily irrelevant targets (i.e. prospective targets during the first search), and targets that

1 were no longer relevant (i.e. previous targets during the second search). This suggests a shared
2 mechanism for preventing interference, whether from future targets or from past targets. The
3 results of Experiment 2 would then imply that items that were never a template for search in the
4 first place, on a certain trial (i.e. the irrelevant condition), would interfere less with the subsequent
5 task and thus did not need to be transformed.

6 An alternative possibility is that local, and arguably more passive mechanisms cause the
7 change in responsiveness. We believe that we can exclude at least two of such mechanisms on the
8 basis of the current data, specifically sensory-induced neural adaptation, and at a macro level, a
9 sensory-induced BOLD undershoot. Both mechanisms would predict the voxels that were most
10 active during memory encoding to be least active a short while later. In fact, such adaptation in
11 responsivity might be functional in memory retrieval (Meyer & Rust, 2018; Turk-Browne, Yi, &
12 Chun, 2006; Ward, Chun, & Kuhl, 2013), and may even also explain the earlier demonstrations of
13 memory reconstruction by Rose et al. (2016) and Wolff et al. (2017). However, Experiment 2 here
14 showed a differential neural response for prospective and irrelevant items, despite the fact that
15 stimulus presentation was identical. Thus, the representational differences we find are determined
16 by the observer's goal state, and not automatically induced by the sensory representation of the
17 stimulus.

18 Although the underlying mechanism remains unknown, we believe the results have
19 important implications for theories of prospective memory storage in working memory. First, the
20 fact that prospectively relevant objects could be decoded from the same regions of interest as the
21 currently relevant objects indicates that the different memory states do not necessarily rely on
22 different brain areas. The successful cross-relevance decoding, where we trained the classifier on
23 one state and tested on another, further confirms this. Second, the idea that prospectively relevant
24 memories are stored in an activity-silent format has recently been debated by Schneegans and Bays
25 (2017) on the basis of the argument that existing data can also be explained by a simpler model
26 which assumes that temporarily irrelevant memories are represented through the same activity as
27 relevant memories, but in a weaker form. Schneegans and Bays (2017) argued specifically against a
28 study by Sprague and colleagues (2016), which indeed showed clear remnants of activity for
29 representations that were assumed to be partly latent. But even when the data reveals no such
30 activity this may only reflect the limited sensitivity of the measure at hand. The reduced activity
31 account is partially supported by our data. We found that during the delay period prior to the first
32 search task, before the evidence for the prospective item diminished to baseline levels, current and
33 prospective representations were highly similar, as evidenced by a strong correlation and

1 successful cross-relevance classification of the current, prospective and irrelevant representations.
2 However, the reduced activity account does not explain that current and prospective
3 representational patterns were very dissimilar during the search. In fact, the partial anti-
4 correlation indicates suppression rather than activation of the relevant voxels.

5 Instead, the emergence of the prospective memory that we found here during the first
6 search fits best with a change in responsivity, resulting in an activity-silent representation. The fact
7 that it was necessary to add activity to the system for the prospective memory to emerge – here in
8 the form of unrelated visual search displays, is already testament to this. Importantly, the current
9 data puts limits on the potential mechanisms by which the responsivity changes. A frequent
10 hypothesis is that prospectively relevant representations are stored through temporary synaptic
11 potentiation. Such short-term potentiation predicts that what was strongly activated during
12 encoding, will become more responsive, when prospective, and fire more strongly when
13 reactivated. We observed the opposite: What was strongly activated when current became more
14 strongly suppressed when prospective and vice versa. This goes against a simple short-term
15 potentiation account of activity-silent representations in working memory.

16 In conclusion, we find evidence that, in trying to separate current and prospective goals in
17 visual search, the brain stores representations within the same neuronal ensembles, but through
18 opposite representational patterns.

19

20 **Methods**

21 **Participants**

22 Twenty-four participants (8 males, $M = 26.74$ years of age, $SD = 3.21$ years) participated in
23 Experiment 1, and twenty-five¹ participants in Experiment 2 (14 males, $M = 25$ years of age, $SD = 4.5$
24 years). For both experiments, we obtained written informed consent from each participant before
25 experimentation. Participants had normal or corrected-to-normal vision. The experiment was
26 approved by the Ethical Committee of the Faculty of Social and Behavioral Sciences, University of
27 Amsterdam (where scanning took place) and conformed to the Declaration of Helsinki.

28 **Task and Stimuli**

¹ Based on Experiment 1, we planned 24 participants. We tested an extra subject to ensure that we would have a complete sample in case a participant had to be excluded. In the end, this turned out unnecessary, and all tested participants were included in the analyses.

1 **Experiment 1.** On each trial, participants performed two consecutive visual search tasks of
2 real-world objects. The object of interest (cow, dresser or skate) consisted of real-world greyscale
3 photographs, selected out of four exemplars. These categories were selected to have maximal
4 dissimilarity in representational space (see Harel et al., 2014). The object of interest (cow, dresser
5 or skate) was to be searched for first or second – thus making it currently or prospectively relevant
6 (referred to as the variable template search). To maximize the chances of decoding the target of
7 interest (whether current or prospective), and to limit the working memory load, the remaining
8 search task always involved the same ‘daisy’ flower target (i.e. Only 1 exemplar) referred to as the
9 constant template search.

10 As can be seen in Figure 1, each trial started with a fixation followed by the sequential
11 presentation of two memory items (variable template [cow, dresser or skate] and constant
12 template [the daisy], 2.4° visual angle) each presented for 750 ms with a 500 ms fixation in
13 between. After a fixation of 500 ms a cue, either a 1 or a 2 was presented indicating the search
14 order in which the memory items needed to be searched for in two subsequent search tasks. Thus,
15 participants either had to search for the variable template first and then the constant daisy
16 template (Current condition) or the daisy had to be searched for first and the variable template
17 second (Prospective condition). Both relevance conditions (Current and Prospective), order of the
18 memory items as well as the cue were counterbalanced across trials. The cue was followed by an 8
19 second delay with a fixation dot in the middle of the screen (‘Delay’) after which the first search
20 display was presented. The search display consisted of 6 different exemplars (2.4° visual angle) of
21 the same category as the cued memory item and could either contain the remembered ‘Current’
22 object (‘Present’) or not (‘Absent’). Participants had to indicate through button presses with their
23 left and right hand whether the memory item was present or absent. The distractors in the search
24 displays were randomly placed among a radius of (7.4° visual angle). The search display was
25 presented for two seconds and participants had to respond within these two seconds. After the
26 first search display another eight seconds blank delay period followed (‘Search 1’) and then the
27 second search display was shown depicting exemplars from the uncued object category. This was
28 again followed by an eight second inter trial interval (ITI) (‘Search 2’). After completion of the first
29 search task, observers had to turn to the prospective item, and indicate its presence or absence in
30 the second search display. At the end of each trial, within the ITI, participants received feedback
31 (for 400 ms) on their performance for each search tasks: either ‘correct’, ‘incorrect’ or ‘missed’ (if
32 the response did not occur within the 2 seconds duration of the search displays). At the end of each

1 run the percentage correct and average reaction times were presented for both the constant
2 template search and the variable template search.

3 Notice that in Experiment 1, the current and prospective templates were always drawn
4 from separate category sets within a trial. Specifically, in current trials, the current item was drawn
5 from one of the three categories in the variable template set (object of interest: cows, dressers or
6 skates) and the prospective item was always the same constant template flower. In prospective
7 trials the category sets reversed, with the prospective item being the variable template. We did this
8 intentionally; having independent sets for the two search templates - within a trial - ensured that
9 we could unequivocally interpret the category classification accuracy as reflecting the
10 representation of the object of interest when either current or prospective. Remember that the
11 classifiers learned to differentiate the neural pattern of the categories of interest: cow, dresser and
12 skate. If both templates were to be drawn from the same category set within the trial, it would be
13 impossible to know whether category classification accuracy actually reflects the quality of the
14 representation of the current template, the prospective template or a combination of the two.

15 The main experiment consisted of 8 runs with 12 trials each (96 trials in total). Each run
16 contained equal amount of trials per condition and category of interest (i.e., cow, dresser and
17 skate). Each experimental run lasted ~ 7 minutes. The total duration of a session was ~1.5 hours
18 (including the structural scan (6 minutes) and mapper run (7 minutes), see below).

19 **Experiment 2.** The task in Experiment 2 was similar to Experiment 1, but we included
20 important changes (see Figure 1B). First, we replaced the constant template search with a
21 duplicate search where participants had to indicate if one of the objects appeared twice in the
22 search display. Second, the duplicate search tasks changed the category from trial to trial to be one
23 out of three possible categories (butterfly, motorcycles and trees). Finally, we added a third
24 condition (i.e., Irrelevant condition), where after the cue participants could immediately drop the
25 item from memory as they only performed the duplicate search task.

26 Each trial started with the presentation of only one memory item (cow, dresser or skate) for
27 1500 ms, followed by a fixation display that stayed on for 1500 ms. Then, a cue indicated the
28 relevance of the memory item. The cue could be either a 1, 2 or 0 and remained on the screen for
29 1000 ms. When the cue was '1', participants performed the template search first and the duplicate
30 search second, making the memorized object currently relevant (Current condition). The order
31 reversed when the cue was '2', rendering the memorized object only prospectively relevant, as

1 observers performed the duplicate search first and the template search second (Prospective
2 condition). Finally, if the cue was '0' the object was irrelevant because participants would only
3 perform the duplicate search (Irrelevant condition) and would not be tested on the memory object.
4 As in Experiment 1, the cue was followed by an 8 second delay with a fixation cross in the middle of
5 the screen ('Delay') after which the first search display was presented. Depending on the condition,
6 the first search display was either a template search (Current condition) or a duplicate search
7 (Prospective and Irrelevant conditions). In the template search, participants indicated with a
8 button press whether the memorized object of interest was present or absent among six exemplars
9 of the same object category. Similarly, in the duplicate search, participants indicated whether a
10 duplicate object (i.e., the same exemplar appeared twice in the search display) was present or
11 absent, again set size for this display was six objects. After the first search display another eight
12 seconds blank delay period followed ('Search 1') after which the trial either ended (Irrelevant
13 condition) or the second search display was presented (Current and Prospective conditions). This
14 second search display was also followed by an eight second blank period ('Search 2') after which
15 the trial ended. The location of items and the duration of the search displays were the same as in
16 Experiment 1.

17 The main experiment consisted of 9 runs with 12 trials each (108 trials in total). Within
18 each experimental run, we balanced the amount of times that each relevance condition (Current,
19 Prospective and irrelevant) was presented (four trials per condition), as well as the amount of
20 times that participants had to respond either 'present' or 'absent' in each search task. However, the
21 relevance condition by category combinations (i.e., 9 in total: 3 relevance conditions [current,
22 prospective, irrelevant] x 3 memory category [cow, dresser, skate]) could not be completely
23 balanced within runs (12 trials per run); nonetheless, across the whole experiment there were
24 equal amount of trials for each combination. We also balanced the category used in the duplicate
25 search task (butterfly, motorcycles and trees) across conditions and in combination with the
26 category of the variable template (i.e. cow, dresser, skate). Each experimental run lasted ~ 7
27 minutes. The total duration of a session was ~1.7 hours (including, short breaks in between runs,
28 the structural scan and mapper run).

29 In both Experiments, the stimuli were back-projected on a 61 x 36 cm LCD screen (1920 x
30 1080 pixels) using Presentation (Neurobehavioral Systems, Albany, CA, USA) and viewed through a
31 mirror attached to the head coil. Eye tracking data (EyeLink 1000, SR Research, Canada) was
32 monitored to ensure participants were awake and attending the stimuli.

1 **Regions of Interest: object-selective cortex mapper (pFs)**

2 At the end of each session we independently mapped the region of interest as the region that
3 responded more strongly to intact vs. scrambled objects (Malach et al., 1995), within an anatomical
4 mask of the temporal occipital fusiform cortex (from the Harvard-Oxford Cortical Structural Atlas of
5 the FSL package). We used the same images and object categories as in our experimental tasks
6 (Experiment 1: cow, skate, dresser and flower; Experiment 2: cow, dresser, skate, butterfly,
7 motorbike and tree). This localized object-selective region of interest corresponded to the posterior
8 fusiform part of lateral occipital cortex (pFs), also referred to as posterior fusiform gyrus (pFG).
9 Stimuli were presented in 24 blocks of images from the same category with each image shown for
10 250 ms. In Experiment 1, stimuli consisted of 48 intact objects (12 of each object category) and 48
11 scrambled objects (12 of each object category) that were presented in separate blocks for each
12 object category (24 blocks in total, 6 per category) with fixation block intermixed (seven in total). In
13 Experiment 2, because we also included the categories from the duplicate search task, we had 72
14 intact objects (12 per category) and 72 scrambled objects also presented in separate blocks for
15 each object category (24 blocks in total, 4 per category).

16 In both Experiments, the mapper run lasted ~ 7 minutes. Participants were asked to push a
17 button when two consecutive images were identical (same exemplar) to ensure attention. The same
18 fMRI preprocessing steps as described for the experimental task were performed for this mapper.
19 In Experiment 1, for two participants the data recorded from this mapper was not usable, therefore
20 we used the anatomical mask only for these participants. In Experiment 2, due to a programming
21 error, for the first 3 participants the mapper only contained three categories (cow, skate and tree).

22

23 **fMRI Acquisition**

24 Scanning was done on a 3T Philips Achieva TX MRI scanner with a 32-elements head coil. In the
25 middle of the testing session (after 4 runs) a high-resolution 3DT1-weighted anatomical image (TR,
26 8.35 ms; TE, 3.83 ms; FOV, $240 \times 220 \times 188$, 1 mm³ voxel size) was recorded for every participant
27 (duration 6 minutes).

28 During the experimental task an object-selective cortex functional localizer, blood oxygenation level
29 dependent (BOLD)-MRI was recorded using Echo Planar Imaging (EPI) (TR 2000 ms, TE 27.62 ms,
30 FA 76.1, 36 slices with ascending acquisition, voxel size 3 mm³, slice gap 0.3 mm, FOV 240 x 118.5 x
31 240).

1 **fMRI Data analysis**

2 **fMRI Preprocessing.**

3 MRI data was registered to the subject specific T1 scan using boundary based registration (Greve &
4 Fischl, 2009). The subject-specific T1 scan was registered to the MNI brain using FMRIB's Nonlinear
5 Image Registration Tool (FNIRT). For the functional imaging data we used FEAT version 5, part of
6 FSL (Oxford Centre for Functional MRI of the Brain (FMRIB) Software Library;
7 www.fmrib.ox.ac.uk/fsl; (S. M. Smith et al., 2004)). Preprocessing steps consisted of motion
8 correction, brain extraction, slice-time correction, alignment, and high-pass filtering (cutoff 100 s).
9 For each subject and each trial a general linear model (GLM) was fitted to the data, whereby every
10 TR (2 seconds each) was taken as a regression variable. We derived the t-value of each voxel for
11 each of the fifteen (Experiment 1) or sixteen (Experiment 2) TRs that were part of each trial in
12 Experiment 1 and Experiment 2 respectively. We used FMRIB's Improved Linear Model (FILM)
13 (Woolrich, Ripley, Brady, & Smith, 2001) for the time-series statistical analysis. The data was further
14 analyzed in Matlab (The MathWorks, Natick, MA, USA). For every participant, every run, every
15 experimental condition (Experiment 1: Current, Prospective; Experiment 2: Current, Prospective
16 and Irrelevant), category exemplar (Cow, Dresser and Skate, 4 exemplars of each) and for each TR,
17 we created a vector containing the t-value per voxel in our regions of interest (see below). T-values
18 for each predictor were computed by dividing the beta-weight by the standard error. That vector
19 comprised the spatial pattern of activity evoked at that time point (TR) for that experimental
20 condition in our region of interest.

21 **Within-relevance and Cross-relevance Object Category Decoding.**

22 Next, we used these multi-voxel patterns to answer the question whether Relevance (Experiment 1:
23 current or prospective; Experiment 2: current, prospective, irrelevant) affected the neural category
24 representations. To determine this we used the Princeton Multi-Voxel Pattern Analysis toolbox
25 (available at <https://github.com/princetonuniversity/princeton-mvpa-toolbox>, see Detre et al.
26 2006). To examine whether current, prospective and irrelevant items evoked a distinct pattern of
27 activity in pFs, for each condition and TR, a single class logistic regression classifier was trained to
28 distinguish each object category (cow, dresser and skate). Logistic regression computes a weighted
29 combination of voxel activity values, and it adjusts the (per-voxel) regression weights to minimize
30 the discrepancy between the predicted output value and the correct output value. The maximum

1 number of iterations used by the iteratively-reweighted least squares (IRLS) algorithm was set to
2 5000.

3 In Experiment 1, classifier performance was evaluated with a standard leave one run out
4 cross validation procedure. This involved training a single class logistic regression classifier to
5 learn a mapping between the neural patterns and the corresponding category labels for all but one
6 run, and then using the trained classifier to predict the category of stimuli from the test patterns in
7 the remaining run. For each iteration we trained the classifier on seven runs and tested on the
8 remaining run for each ROI. Overall classification accuracy was the average accuracy of the nine
9 iterations.

10 In Experiment 2, relevance condition was fully balanced within each run; however, we could
11 not fully balance the relevance condition by object category combinations within each run (see Task
12 and stimuli of Experiment 2). Therefore, we used a modified leave one run out cross-validation
13 procedure. Per run we had four trials per relevance condition; therefore, each training set
14 consisted of 8 runs with 32 trials per relevance condition which is not a multiple of 3 (i.e., the
15 amount object categories of interest). Thus, for each relevance condition (Current, Prospective
16 Irrelevant) when selecting all but one run for the training set, one of the categories contained 10
17 exemplars while the other two categories had 11. Likewise, the testing set contained all three
18 categories (i.e., cow, dresser, skate), but two of the four exemplars belonged to the category that
19 was less frequent in the training set. To correct for this slight unbalance and ensure that the
20 classifiers were not biased against the least frequent category, we picked one exemplar (for each of
21 the two categories that had 11) and excluded them from the training set, leaving 30 exemplars per
22 training set per relevance condition (10 per category). We repeated this entire process and left a
23 different exemplar out of the training set, until all exemplars were left out exactly once and all
24 exemplars were included an equal number of times across all train-test procedures. Therefore, for
25 each run out, we trained and tested 11 classifiers (99 in total per TR). Overall classification
26 accuracy was the average accuracy of these 99 iterations. Moreover, we used a balanced accuracy
27 calculation as described in Fahrenfort et al. (2018), where accuracy is calculated separated per
28 class and then averaged across classes.

29 In both Experiments we ran two types of decoding. We investigated category decoding
30 (Cow, Dresser and Skate) both *within* Current, Prospective and Irrelevant conditions (within-
31 relevance decoding) and *between* relevance conditions (cross-relevance decoding) for each time
32 point (TR) in the trial separately. For the within relevance classification, we trained and tested the

1 classifier on the same condition (Current, Prospective or Irrelevant). For the cross-relevance
2 classification in Experiment 1, we trained when the category was a Current item and tested when
3 the category was a Prospective item ('PC') and vice versa ('CP'). In Experiment 2, we applied this
4 same cross-relevance decoding scheme (Current-Prospective) and added two more: Current-
5 Irrelevant and Prospective-Irrelevant. This resulted in six different testing and training
6 combinations. To reduce the amount of comparisons needed, we averaged the classification
7 performance of those combinations where the same conditions were used either for testing or
8 training. All the significance testing was performed on the averaged data of Current-Prospective,
9 Current-Irrelevant and Prospective-Irrelevant. We obtained a classification score (percentage
10 correct) per participant for every relevance condition and time point (TR). Note that here chance
11 decoding was 33.33% since we had three object categories (Cow, Dresser and Skate). All statistical
12 comparisons are based on two-tailed tests, except for the comparison against chance in the within-
13 relevance coding scheme as there decoding cannot go below chance (cf. Christophel et al., 2018). All
14 statistical analyses were performed using SPSS 17.0 (IBM, Armonk, USA).

15 **Representational dissimilarity Analysis.**

16 For each TR we created a representational dissimilarity matrix (RDM) (Kriegeskorte et al., 2008;
17 Kriegeskorte & Kievit, 2013). Each cell of the matrix represents a 1-rho (Spearman correlation) of
18 the activity patterns of two individual exemplars. In Experiment 1, the RDMs consisted of 24x24, 4
19 unique exemplars per category (Cow, Dresser and Skate) and 2 different relevance levels (Current
20 and Prospective). In Experiment 2, we created three different sets of RDMs depending on the
21 conditions correlated (Current-Prospective, Current-Irrelevant, Prospective-Irrelevant). The RDMs
22 of each run were averaged to obtain one RDM per TR. We further averaged across the three TRs for
23 each interval of interest in the trial (Delay, Search 1 and Search 2). For visualization purposes we
24 transformed the RDM by replacing each element by its rank in the distribution of all its elements
25 (scaled between 0 to 1). In addition, we used multidimensional scaling (MDS) plots wherein the
26 distance between points reflects the dissimilarity in their neural patterns of response. To compute
27 the interaction between Relevance and Category over the course of the trial we calculated the
28 dissimilarity for the between Relevance (Experiment 1: Current-Prospective, Experiment 2: Current
29 -Prospective; Current-Irrelevant; Prospective-Irrelevant) and Category (same [Cow/Cow,
30 Dresser/Dresser and Skate/Skate] vs different [Cow/Dresser, Dresser/Skate, Skate/Dresser]) by
31 averaging the cells within each class. We calculated this for every TR separately, and then averaged
32 those across the three TRs in the predetermined intervals (Delay, Search 1 and Search 2).

1

References

- 2 Albers, A. M., Kok, P., Toni, I., Dijkerman, H. C., & de Lange, F. P. (2013). Shared Representations for
3 Working Memory and Mental Imagery in Early Visual Cortex. *Current Biology*, 23(15), 1427–
4 1431. <http://doi.org/10.1016/j.cub.2013.05.065>
- 5 Anderson, M. C., Ochsner, K. N., Kuhl, B., Cooper, J., Robertson, E., Gabrieli, S. W., et al. (2004). Neural
6 systems underlying the suppression of unwanted memories. *Science*, 303(5655), 232–235.
7 <http://doi.org/10.1126/science.1089504>
- 8 Banich, M. T., Mackiewicz Seghete, K. L., Depue, B. E., & Burgess, G. C. (2015). Multiple modes of
9 clearing one's mind of current thoughts: overlapping and distinct neural systems.
10 *Neuropsychologia*, 69, 105–117. <http://doi.org/10.1016/j.neuropsychologia.2015.01.039>
- 11 Barak, O., & Tsodyks, M. (2014). Working models of working memory. *Current Opinion in*
12 *Neurobiology*, 25, 20–24. <http://doi.org/10.1016/j.conb.2013.10.008>
- 13 Carlisle, N. B., & Woodman, G. F. (2011). Automatic and strategic effects in the guidance of attention
14 by working memory representations. *Acta Psychologica*, 137(2), 217–225.
15 <http://doi.org/10.1016/j.actpsy.2010.06.012>
- 16 Christophel, T. B., Iamshchinina, P., Yan, C., Allefeld, C., & Haynes, J.-D. (2018). Cortical specialization
17 for attended versus unattended working memory. *Nature Neuroscience*.
18 <http://doi.org/10.1038/s41593-018-0094-4>
- 19 Çukur, T., Nishimoto, S., Huth, A. G., & Gallant, J. L. (2013). Attention during natural vision warps
20 semantic representation across the human brain. *Nature Neuroscience*, 16(6), 763–770.
21 <http://doi.org/10.1038/nn.3381>
- 22 de Vries, I. E. J., van Driel, J., Karacaoglu, M., & Olivers, C. N. L. (2018). Priority Switches in Visual
23 Working Memory are Supported by Frontal Delta and Posterior Alpha Interactions. *Cerebral*
24 *Cortex*. <http://doi.org/10.1093/cercor/bhy223>
- 25 Depue, B. E., Curran, T., & Banich, M. T. (2007). Prefrontal regions orchestrate suppression of
26 emotional memories via a two-phase process. *Science*, 317(5835), 215–219.
27 <http://doi.org/10.1126/science.1139560>
- 28 D'Esposito, M., & Postle, B. R. (2015). The cognitive neuroscience of working memory. *Annual*
29 *Review of Psychology*, 66(1), 115–142. <http://doi.org/10.1146/annurev-psych-010814-015031>
- 30
- 31 Detre, G., Polyn, S. M., Moore, C. D., Natu, V. S., Singer, B., Cohen, J. D., ... & Norman, K. A. (2006, June).
32 The multi-voxel pattern analysis (MVPA) toolbox. In Poster presented at the annual Meeting of
33 the organization for human brain mapping.
- 34 Fahrenfort JJ, van Driel J, van Gaal S, Olivers CNL (2018) From ERPs to MVPA using the
35 Amsterdam Decoding and Modeling Toolbox (ADAM). *Front Neurosci*, 12(368),
36 <https://doi.org/10.3389/fnins.2018.00368>
- 37 Downing, P. E., & Dodds, C. M. (2003). Competition in visual working memory for control of search.
38 *Perception*, 32, 92–93.
- 39 Erickson, M. A., Maramba, L. A., & Lisman, J. (2010). A Single Brief Burst Induces GluR1-dependent
40 Associative Short-term Potentiation: A Potential Mechanism for Short-term Memory. *Journal of*
41 *Cognitive Neuroscience*, 22(11), 2530–2540. <http://doi.org/10.1162/jocn.2009.21375>
- 42 Greve, D. N., & Fischl, B. (2009). Accurate and robust brain image alignment using boundary-based
43 registration. *Neuroimage*, 48(1), 63–72. <http://doi.org/10.1016/j.neuroimage.2009.06.060>
- 44 Harel, A., Kravitz, D. J., & Baker, C. I. (2014). Task context impacts visual object processing
45 differentially across the cortex. *Proceedings of the National Academy of Sciences of the United*
46 *States of America*, 111(10), E962–71. <http://doi.org/10.1073/pnas.1312567111>
- 47 Harrison, S. A., & Tong, F. (2009). Decoding reveals the contents of visual working memory in early
48 visual areas. *Nature*, 458(7238), 632–635. <http://doi.org/10.1038/nature07832>

- 1 Henson, R., & Rugg, M. D. (2003). Neural response suppression, haemodynamic repetition effects,
2 and behavioural priming. *Neuropsychologia*, 41(3), 263–270.
- 3 Houtkamp, R., & Roelfsema, P. R. (2006). The effect of items in working memory on the deployment
4 of attention and the eyes during visual search. *Journal of Experimental Psychology: Human*
5 *Perception and Performance*, 32(2), 423–442. <http://doi.org/10.1037/0096-1523.32.2.423>
- 6 Huettel, S. A., & McCarthy, G. (2000). Evidence for a refractory period in the hemodynamic response
7 to visual stimuli as measured by MRI. *Neuroimage*, 11(5), 547–553.
8 <http://doi.org/10.1006/nimg.2000.0553>
- 9 King, J.-R., & Dehaene, S. (2014). Characterizing the dynamics of mental representations: the
10 temporal generalization method. *Trends in Cognitive Sciences*, 18(4), 203–210.
11 <http://doi.org/10.1016/j.tics.2014.01.002>
- 12 Kravitz, D. J., Saleem, K. S., Baker, C. I., Ungerleider, L. G., & Mishkin, M. (2013). The ventral visual
13 pathway: an expanded neural framework for the processing of object quality. *Trends in*
14 *Cognitive Sciences*, 17(1), 26–49. <http://doi.org/10.1016/j.tics.2012.10.011>
- 15 Kriegeskorte, N., & Kievit, R. A. (2013). Representational geometry: integrating cognition,
16 computation, and the brain. *Trends in Cognitive Sciences*, 17(8), 401–412.
17 <http://doi.org/10.1016/j.tics.2013.06.007>
- 18 Kriegeskorte, N., Mur, M., Ruff, D. A., Kiani, R., Bodurka, J., Esteky, H., et al. (2008). Matching
19 Categorical Object Representations in Inferior Temporal Cortex of Man and Monkey. *Neuron*,
20 60(6), 1126–1141. <http://doi.org/10.1016/j.neuron.2008.10.043>
- 21 LaRocque, J. J., Lewis-Peacock, J. A., & Postle, B. R. (2014). Multiple neural states of representation in
22 short-term memory? It's a matter of attention. *Frontiers in Human Neuroscience*, 8.
23 <http://doi.org/10.3389/fnhum.2014.00005>
- 24 LaRocque, J. J., Lewis-Peacock, J. A., Drysdale, A. T., Oberauer, K., & Postle, B. R. (2013). Decoding
25 Attended Information in Short-term Memory: An EEG Study. *Journal of Cognitive Neuroscience*,
26 25(1), 127–142. http://doi.org/10.1162/jocn_a_00305
- 27 LaRocque, J. J., Riggall, A. C., Emrich, S. M., & Postle, B. R. (2017). Within-Category Decoding of
28 Information in Different Attentional States in Short-Term Memory. *Cerebral Cortex*, 27(10),
29 4881–4890. <http://doi.org/10.1093/cercor/bhw283>
- 30 Larsson, J., & Smith, A. T. (2012). fMRI Repetition Suppression: Neuronal Adaptation or Stimulus
31 Expectation? *Cerebral Cortex*, 22(3), 567–576. <http://doi.org/10.1093/cercor/bhr119>
- 32 Lee, S.-H., Kravitz, D. J., & Baker, C. I. (2013). Goal-dependent dissociation of visual and prefrontal
33 cortices during working memory. *Nature Neuroscience*, 16(8), 997–999.
34 <http://doi.org/10.1038/nn.3452>
- 35 Lewis-Peacock, J. A., & Postle, B. R. (2012). Decoding the internal focus of attention.
36 *Neuropsychologia*, 50(4), 470–478. <http://doi.org/10.1016/j.neuropsychologia.2011.11.006>
- 37 Malach, R., Reppas, J. B., Benson, R. R., Kwong, K. K., Jiang, H., Kennedy, W. A., et al. (1995). Object-
38 Related Activity Revealed by Functional Magnetic Resonance Imaging in Human Occipital
39 Cortex. *Proceedings of the National Academy of Sciences of the United States of America*,
40 92(18), 8135–8139. [http://doi.org/10.2307/2368011?ref=search-](http://doi.org/10.2307/2368011?ref=search-gateway:ae8ad9f1e027b5b5af79d6209b7141b9)
41 [gateway:ae8ad9f1e027b5b5af79d6209b7141b9](http://doi.org/10.2307/2368011?ref=search-gateway:ae8ad9f1e027b5b5af79d6209b7141b9)
- 42 Mallett, R., & Lewis-Peacock, J. A. (2018). Behavioral decoding of working memory items inside and
43 outside the focus of attention. *Annals of the New York Academy of Sciences*, 12, 342.
44 http://doi.org/10.1162/jocn_a_01180
- 45 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data. *Journal*
46 *of Neuroscience Methods*, 164(1), 177–190. <http://doi.org/10.1016/j.jneumeth.2007.03.024>
- 47 Meyer, T., & Rust, N. (2018). Single-exposure visual memory judgments are reflected in
48 inferotemporal cortex. *eLife*, 7. <http://doi.org/10.7554/eLife.32259>
- 49 Mongillo, G., Barak, O., & Tsodyks, M. (2008). Synaptic theory of working memory. *Science*,
50 319(5869), 1543–1546. <http://doi.org/10.1126/science.1150769>

- 1 Myers, N. E., Stokes, M. G., & Nobre, A. C. (2017). Prioritizing Information during Working Memory:
2 Beyond Sustained Internal Attention. *Trends in Cognitive Sciences*, 21(6), 449–461.
3 <http://doi.org/10.1016/j.tics.2017.03.010>
- 4 Nastase, S. A., Connolly, A. C., Oosterhof, N. N., Halchenko, Y. O., Guntupalli, J. S., Castello, M. V. D. O.,
5 et al. (2017). Attention Selectively Reshapes the Geometry of Distributed Semantic
6 Representation. *Cerebral Cortex*, 27(8), 4277–4291. <http://doi.org/10.1093/cercor/bhx138>
7 Olivers & Eimer, 2011
- 8 Olivers, C. N. L., & Eimer, M. (2011). On the difference between working memory and attentional set.
9 *Neuropsychologia*, 49(6), 1553–1558.
10 <http://doi.org/10.1016/j.neuropsychologia.2010.11.033>
- 11 Olivers, C. N. L., Peters, J., Houtkamp, R., & Roelfsema, P. R. (2011). Different states in visual working
12 memory: When it guides attention and when it does not. *Trends in Cognitive Sciences*, 15(7),
13 327–334.
- 14 Peters, J. C., Roelfsema, P. R., & Goebel, R. (2012). Task-Relevant and Accessory Items in Working
15 Memory Have Opposite Effects on Activity in Extrastriate Cortex. *The Journal of Neuroscience*,
16 32(47), 17003–17011. <http://doi.org/10.1523/JNEUROSCI.0591-12.2012>
- 17 Reddy, L., Kanwisher, N. G., & VanRullen, R. (2009). Attention and biased competition in multi-voxel
18 object representations. *Proceedings of the National Academy of Sciences of the United States of*
19 *America*, 106(50), 21447–21452. <http://doi.org/10.1073/pnas.0907330106>
- 20 Reeder, R. R., Olivers, C. N. L., & Pollmann, S. (2017). Cortical evidence for negative search
21 templates. *Visual Cognition*, 25(1-3), 278–290. [http://doi.org/10.1016/0306-4522\(83\)90265-](http://doi.org/10.1016/0306-4522(83)90265-8)
22 8
- 23 Rose, N. S., LaRocque, J. J., Riggall, A. C., Gosseries, O., Starrett, M. J., Meyering, E. E., & Postle, B. R.
24 (2016). Reactivation of latent working memories with transcranial magnetic stimulation.
25 *Science*, 354(6316), 1136–1139. <http://doi.org/10.1126/science.aah7011>
- 26 Schneegans, S., & Bays, P. M. (2017). Restoration of fMRI Decodability Does Not Imply Latent
27 Working Memory States. *Journal of Cognitive Neuroscience*, 29(12), 1977–1994.
28 <http://doi.org/10.1038/nature04262>
- 29 Serences, J. T., Ester, E. F., Vogel, E. K., & Awh, E. (2009). Stimulus-Specific Delay Activity in Human
30 Primary Visual Cortex. *Psychological Science*, 20(2), 207–214. [http://doi.org/10.1111/j.1467-](http://doi.org/10.1111/j.1467-9280.2009.02276.x)
31 9280.2009.02276.x
- 32 Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E. J., Johansen-Berg, H., et
33 al. (2004). Advances in functional and structural MR image analysis and implementation as FSL.
34 *Neuroimage*, 23, S208–S219. <http://doi.org/10.1016/j.neuroimage.2004.07.051>
- 35 Sprague, T. C., Ester, E. F., & Serences, J. T. (2016). Restoring Latent Visual Working Memory
36 Representations in Human Cortex. *Neuron*, 91(3), 694–707.
37 <http://doi.org/10.1016/j.neuron.2016.07.006>
- 38 Stokes, M. G. (2015). “Activity-silent” working memory in prefrontal cortex: a dynamic coding
39 framework. *Trends in Cognitive Sciences*, 19(7), 394–405.
40 <http://doi.org/10.1016/j.tics.2015.05.004>
- 41 Sugase-Miyamoto, Y., Liu, Z., Wiener, M. C., Optican, L. M., & Richmond, B. J. (2008). Short-Term
42 Memory Trace in Rapidly Adapting Synapses of Inferior Temporal Cortex. *PLoS Computational*
43 *Biology*, 4(5), e1000073. <http://doi.org/10.1371/journal.pcbi.1000073.g012>
- 44 Turk-Browne, N. B., Yi, D.-J., & Chun, M. M. (2006). Linking Implicit and Explicit Memory: Common
45 Encoding Factors and Shared Representations. *Neuron*, 49(6), 917–927.
46 <http://doi.org/10.1016/j.neuron.2006.01.030>
- 47 van Loon, A. M., Olmos Solis, K., & Olivers, C. N. L. (2017). Subtle eye movement metrics reveal task-
48 relevant representations prior to visual search. *Journal of Vision*, 17(6), 13,
49 <http://doi.org/10.1167/17.6.13>

- 1
2 Vautin, & Berkley, M. A. (1977). Responses of Single Cells in Cat Visual-Cortex to Prolonged
3 Stimulus Movement - Neural Correlates of Visual Aftereffects. *Journal of Neurophysiology*,
4 40(5), 1051–1065. <http://doi.org/10.1152/jn.1977.40.5.1051>
5 Ward, E. J., Chun, M. M., & Kuhl, B. A. (2013). Repetition suppression and multi-voxel pattern
6 similarity differentially track implicit and explicit visual memory. *The Journal of Neuroscience*,
7 33(37), 14749–14757. <http://doi.org/10.1523/JNEUROSCI.4889-12.2013>
8 Wimber, M., Alink, A., Charest, I., Kriegeskorte, N., & Anderson, M. C. (2015). Retrieval induces
9 adaptive forgetting of competing memories via cortical pattern suppression. *Nature*
10 *Neuroscience*, 18(4), 582–589. <http://doi.org/10.1038/nn.3973>
11 Wolff, M. J., Ding, J., Myers, N. E., & Stokes, M. G. (2015). Revealing hidden states in visual working
12 memory using electroencephalography. *Frontiers in Systems Neuroscience*, 9(5), 1427.
13 <http://doi.org/10.1146/annurev.physiol.64.092501.114547>
14 Wolff, M. J., Jochim, J., Akyurek, E. G., & Stokes, M. G. (2017). Dynamic hidden states underlying
15 working-memory-guided behavior. *Nature Neuroscience*, 20(6), 864–.
16 <http://doi.org/10.1038/nn.4546>
17 Woolrich, M. W., Ripley, B. D., Brady, M., & Smith, S. M. (2001). Temporal autocorrelation in
18 univariate linear modeling of FMRI data. *Neuroimage*, 14(6), 1370–1386.
19 <http://doi.org/10.1006/nimg.2001.0931>
20 Yu, Q., & Postle, B. R. (n.d.). Different states of priority recruit different neural codes in visual
21 working memory. <http://doi.org/10.1101/334920>
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

1 Supplementary material

2 Behavioral results

3 **Experiment 1.** We computed mean RTs and mean accuracy for both types of search task
 4 (variable template search and constant template search) when they came either first or second in
 5 the trial. These measures were each entered in a two-way repeated measures ANOVA (N= 24) with
 6 factors search order (Search 1 and Search 2) and type of search task (constant versus variable
 7 template). As expected, the constant template search was overall faster (RT: $F_{(1,23)} = 928.18, p <$
 8 $0.001, \eta_p^2 = 0.97$) and more accurate (percentage correct: $F_{(1,23)} = 183.06, p < 0.001, \eta_p^2 = 0.88$) than
 9 the variable template search. Furthermore, the first search was more accurate than the second
 10 ($F_{(1,23)} = 10.87, p = 0.003$). There was also an interaction for both accuracy and speed ($F_{(1,23)} = 9.22,$
 11 $p = 0.006, \eta_p^2 = 0.28$ and $F_{(1,23)} = 6.32, p = 0.02, \eta_p^2 = 0.21$ respectively): the variable template search
 12 had the lowest accuracy when performed second, while the constant template search was fastest
 13 when second. Overall these results show that we were successful at minimizing the working
 14 memory load for the constant template (flower), thus maximizing the chances of decoding the
 15 variable target category, which was the target of interest.

16 **Figure 1-table Supplement 1.** Percentage correct and RT for Current and Prospective conditions
 17 in Search 1 and Search 2 (N=24) as a function of search order.

	Current		Prospective	
	Search 1	Search 2	Search 1	Search 2
Template	Variable	Constant	Variable	Constant
P. Correct (%)	82.2 (7.1)	98.1 (2.0)	98.6 (2.3)	76.0 (9.9)
RT (ms)	1387(20)	772 (21)	794 (22)	1411 (22)

18
 19
 20 **Experiment 2.** We computed mean RTs and mean percentage correct for the template
 21 search and the duplicate search when they were presented either first or second in the trial. Figure
 22 1-table Supplement 2 shows the mean Percentage correct and RT for Search 1 and Search 2 for each
 23 relevance condition. Next, we ran a two-way repeated measures ANOVA (N= 25) with factors
 24 search order (Search 1 and Search 2) and type of task (template search and duplicate search), using
 25 data from the Current and Prospective conditions only (as the Irrelevant condition only had one

1 search). There were no differences in accuracy between the two types of task as neither the main
 2 effects of search order ($F_{(1,24)} = 1.65, p = 0.210, \eta_p^2 = 0.065$), type of task ($F_{(1,24)} = 0.16, p = 0.747, \eta_p^2$
 3 $= 0.004$), nor their interaction ($F_{(1,24)} = 2.48, p = 0.128, \eta_p^2 = 0.094$) were significant. However, overall
 4 participants were faster in the first than in the second search task ($F_{(1,24)} = 37.24, p < 0.001, \eta_p^2 =$
 5 0.60) and faster in the template search than in the duplicate search ($F_{(1,24)} = 14.76, p < 0.001, \eta_p^2 =$
 6 0.38). There was also a significant interaction between the two factors ($F_{(1,24)} = 19.56, p < 0.001, \eta_p^2$
 7 $= 0.44$), reflecting the fact that participants were faster in the template search when it occurred first
 8 (i.e. Current condition) than second (i.e. Prospective condition, $t_{(1,24)} = -7.40, p < 0.001, d = -1.48$),
 9 while for the duplicate search, speed was the same regardless of the order ($t_{(1,24)} = -1.47, p = .154,$
 10 $d = -0.294$).

11 So while in Experiment 1 the search of interest (variable template) was more difficult than
 12 the remaining search task (constant template), here, if anything, it was the other way around: The
 13 search of interest (template search) was slightly easier than the remaining task (duplicate search).
 14 Moreover, participants were equally accurate at finding the template object regardless of the search
 15 order, but they were slower in the template search when performed second (i.e., Prospective
 16 condition), which suggests that the quality of the memory representation did not decay when it had
 17 to be postponed, but reactivating it for Search 2 required additional time. So any differences
 18 between current and prospective representations were not simply due to participants being worse
 19 on the prospective memory.

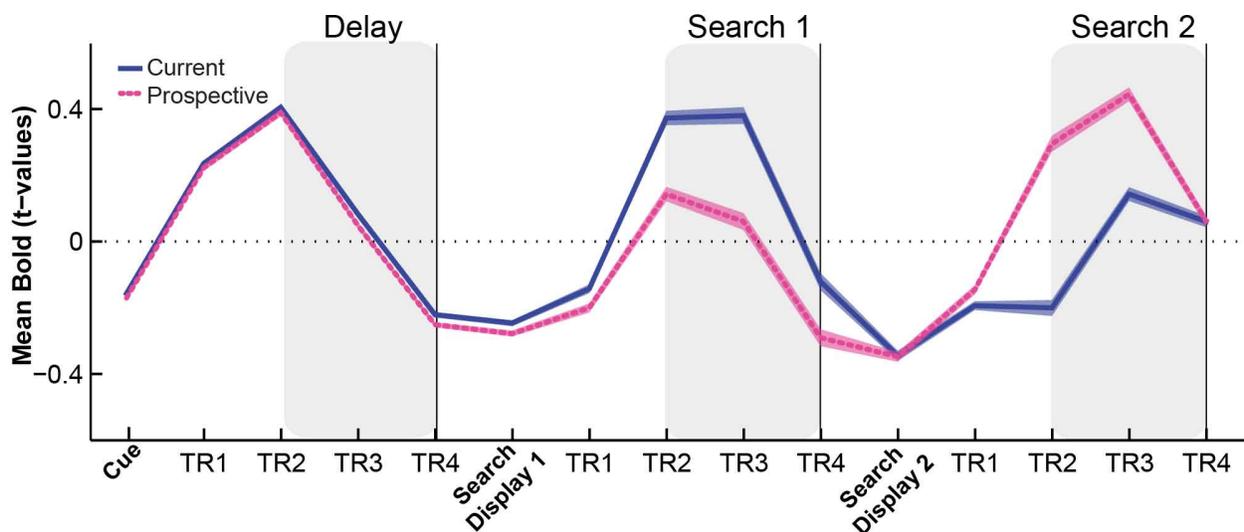
20
 21 **Figure 1-table Supplement 2.** Percentage correct and RT in Search 1 and Search 2 (N=25) as a
 22 function of condition.

	Current		Prospective		Irrelevant
	Search 1	Search 2	Search 1	Search 2	Search 1
	template	Duplicate	template	Duplicate	Duplicate
P. Correct (%)	86.0 (8.3)	84.3 (6.9)	83.9 (5.5)	83.2 (8.4)	83.4 (6.9)
RT (ms)	1355 (96)	1478 (114)	1460 (90)	1147(113)	1469 (91)

23
 24 **Mean BOLD response Experiment 1**

25 Figure 2- figure supplement 1 shows the mean BOLD response in area pFs of Experiment 1. For all

1 the relevance conditions, averaged for the same individual ROIs as used in the MVPA. As can be
2 seen in the figures, for Experiment 1, there was a small difference in the BOLD response magnitude
3 during the Delay depending on whether the category was currently or prospectively relevant ($t_{(1,23)}$
4 = 2.15, $p = 0.0427$, $\eta_p^2 = 0.44$). A stronger difference became apparent for Search 1 ($t_{(1,23)} = 14.46$, p
5 < 0.001, $\eta_p^2 = 1.77$) and Search 2 ($t_{(1,23)} = -13.08$, $p < 0.001$, $\eta_p^2 = -2.67$), where variable template
6 search displays (containing the object categories of interest) elicited a higher response than
7 constant template search displays (with the repeated flower target), probably because the latter
8 was an easier task. .
9



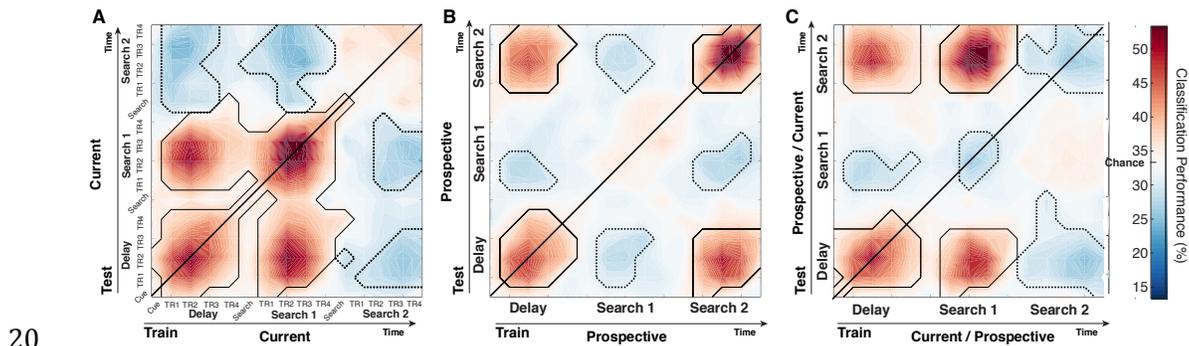
10
11 **Figure 2 - figure Supplement 1:** Time course of the Mean BOLD response Experiment 1. Shaded areas
12 indicate within-subjects s.e.m.
13
14

15 **Cross-temporal generalization of object category decoding in Experiment 1**

16 The main analyses of Experiment 1 are based on decoding performance where training and testing
17 occurred separately for each TR. To examine whether the neural representations of the different
18 object categories for current and prospective states were also related across time, we tested for
19 cross-temporal generalization of decoding accuracy (see King & Dehaene, 2014), by training the
20 classifier on each of the TRs and then testing it on all other TRs in the trial. This was then repeated
21 for all TRs, creating a two-dimensional matrix of cross-temporal object category decoding (with no
22 additional smoothing). Time windows of significant decoding were identified using 2-dimensional
23 cluster-based permutation testing (i.e., across both time axes) with cluster correction ($p = 0.05$ and

1 10.000 iterations) to statistically compare the object category decoding with chance (33.33%)
2 (Maris & Oostenveld, 2007) using and Matlab (The MathWorks, Natick, MA, USA). As a result, we
3 were able to assess the temporal stability of object category decoding and to test whether encoding
4 and maintenance of the object (Delay) was similar to searching for an object (Search 1 and Search
5 2).

6 Figure 2- figure supplement 2 shows the resulting temporal generalization matrices, with
7 red indicating above-chance decoding performance and blue indicating below-chance decoding
8 performance. Note that the pattern on the diagonal reflects the classification per TR as in Figure 2.
9 Of additional interest here are the significant off-diagonal clusters as they indicate a generalized
10 representation across time. These clusters show that target representations maintained during
11 maintenance prior to the first search were similar to the pattern of activity that is observed when
12 subjects are looking at these object categories during search (i.e. the first search when current; see
13 the red off-diagonal clusters in Figure 2- figure supplement 2A; and the second search when
14 prospective; see the red off-diagonal clusters in Figure 2- figure supplement 2B). Figure 2- figure
15 supplement 2C shows the temporal generalization for the cross-relevance decoding scheme
16 (averaged across PC and CP schemes). Notable here is that the same target representations as
17 maintained during maintenance prior to the first search are dissimilar to the prospective
18 representations during the first search, and the no longer relevant representations during the
19 second search, as indicated by the blue off-diagonal clusters.



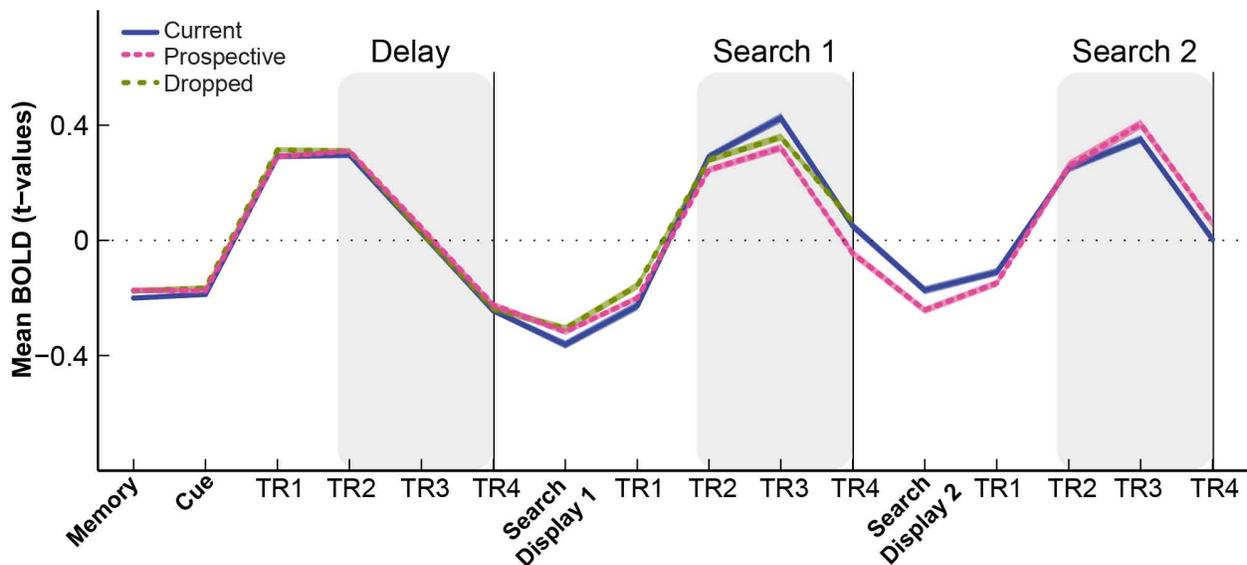
21 **Figure 2- figure supplement 2. Cross-temporal generalization matrices for object category**
22 **decoding as a function of Relevance, with (A) the Current condition, (B) the Prospective condition,**
23 **and (C) cross-relevance classification (averaged for training on current and testing on prospective**
24 **with the transpose of the reverse training scheme). We observed that target representations for the**
25 **objects of interest maintained during the Delay period prior to the first search generalized to the**
26 **search displays that contained that target (i.e. the first search when current, or the second search**
27 **when prospective), as shown by reliable off-diagonal red clusters. In contrast, off-diagonal blue**

1 clusters emerge when the object trained during the delay is prospective during Search 1, or no longer
2 relevant during Search 2, indicating anti-correlated representations. Outlines indicate cluster-based
3 permutation tests with $p < .05$, positive clusters (solid lines), negative clusters (dotted lines), $N=24$.

4 Mean BOLD response Experiment 2

5 Figure 4-figure supplement 1 shows the mean BOLD response in area pFs of Experiment 2. In
6 Experiment 2, there were not significant differences across conditions during the Delay ($F_{(2,48)} =$
7 0.80 , $p = 0.453$, $\eta_p^2 = 0.032$, see Figure S1C and S1D). There was a reliable main effect of condition
8 during the Search 1 interval ($F_{(2,48)} = 13.57$, $p < 0.001$, $\eta_p^2 = 0.361$), with the prospective condition
9 showing a somewhat weaker response than the current ($t_{(1,24)} = 4.32$, $p < 0.001$, $d = 0.86$) and
10 irrelevant ($t_{(1,24)} = 4.66$, $p < 0.001$, $d = 0.93$) conditions, although not as pronounced as in
11 Experiment 1. The Current and Irrelevant condition did not differ from each other ($t_{(1,24)} = 1.30$, $p =$
12 0.205 , $d = 0.26$). As in Experiment 1, the pattern reversed for the Search 2 interval, with the
13 prospective condition showing a stronger BOLD response than the current condition ($t_{(1,24)} = -2.48$,
14 $p = 0.020$, $d = -0.496$), although again less pronounced than in Experiment 1.

15



16

17 **Figure 4 - figure Supplement 1:** Time course of the Mean BOLD response Experiment 2. Shaded areas
18 indicate within-subjects s.e.m.

19

20