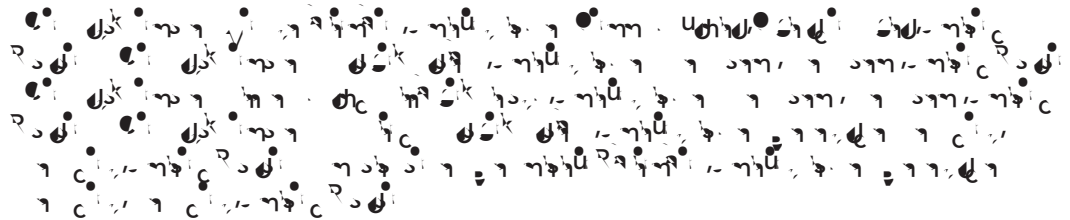


Strategy-dependent effects of working-memory limitations on human perceptual decision-making

Kyra Schapiro^{1*}, Krešimir Josić^{2,3}, Zachary P Kilpatrick^{4,5}, Joshua I Gold¹



Abstract Deliberative decisions based on an accumulation of evidence over time depend on working memory, and working memory has limitations, but how these limitations affect deliberative decision-making is not understood. We used human psychophysics to assess the impact of working-memory limitations on the fidelity of a continuous decision variable. Participants decided the average location of multiple visual targets. This computed, continuous decision variable degraded with time and capacity in a manner that depended critically on the strategy used to form the decision variable. This dependence reflected whether the decision variable was computed either: (1) immediately upon observing the evidence, and thus stored as a single value in memory; or (2) at the time of the report, and thus stored as multiple values in memory. These results provide important constraints on how the brain computes and maintains temporally dynamic decision variables.

Editor's evaluation

This paper employs sophisticated modeling of human behavior in well-controlled tasks to study how limitations of working memory constrain decision-making. Because both are key cognitive processes, that have so far largely been studied in isolation, the paper will be of broad interest to neuroscientists and psychologists. The observed working memory limitations support previous findings and extend them in critical ways.

Introduction

Many perceptual, memory-based, and reward-based decisions depend on an accumulation of evidence over time (*Brody and Hanks, 2016; Gold and Shadlen, 2007; Ratcliff et al., 2016; Shadlen and Shohamy, 2016; Summerfield and Tsetsos, 2012*). This dynamic process, which can operate on timescales ranging from tens to hundreds of milliseconds for many perceptual decisions to seconds or longer for reward-based and other decisions (*Bernacchia et al., 2011; Gold and Stocker, 2017*), requires working memory to maintain representations of new, incoming evidence and/or the aggregated, updating decision variable. Working memory is constrained by capacity and temporal limita-

For spatial working-memory tasks, the precision of working memory for perceived spatial locations is often well described by diffusion dynamics (*Compte et al., 2000; Kilpatrick, 2018; Kilpatrick et al., 2013; Laing and Chow, 2001*) that are commonly implemented in ‘bump-attractor’ models of working memory (*Compte et al., 2000; Constantinidis et al., 2018; Laing and Chow, 2001; Riley and Constantinidis, 2016; Wei et al., 2012; Wimmer et al., 2014*). Our analyses built on this framework by examining memory diffusion dynamics for the different task conditions and potential decision strategies. For the conditions we tested, most participants’ behavior was well fit by one of two distinct strategies, each with its own constraints on decision performance based on different working-memory demands. The first strategy was to compute the decision variable (mean disk angle) immediately upon observing the evidence (individual disk angles), and then store that value in working memory in a manner that, like for the memory of a single perceived angle, could be modeled as a single particle with a particular diffusion constant (Average-then-Diffuse model; AtD). The second strategy was to maintain representations of all disk locations in working memory, modeled as separate diffusing particles, and then to combine them into a decision variable only at the time of the decision (Diffuse-then-Average model; DtA). Such a strategy results in an effective diffusion constant for the average that is inversely related to the number of items. Our results show that like perceived locations, memory for computed mean locations degraded with increased set size (of relevant information), and delay between presentation and report. However, the degree of degradation depended on the strategy used to compute the decision variables, implying that multip utinct lb0.5ai. On P

we measured the error between reported and probed angles as a proxy for working-memory representations and inferred rates of memory degradation (diffusion constants) from the increase in variance of these errors over time within a framework of diffusing-particle models. Below we first describe the model framework, detailing its key assumptions and predictions. We next describe results from Simultaneous conditions, in which all items were presented simultaneously at the beginning of each trial, which demonstrate how capacity and temporal constraints on working memory relate to the accuracy of computed decision variables. We then describe results from Sequential conditions, in which one item was presented after the others in each trial, which demonstrate how capacity and temporal constraints affect the process of evidence integration over time.

Diffusing-particle framework and predictions

Within our diffusing-particle framework, the memory of an item is represented by the location of a diffusing particle. This representation allows us to quantify the corruption (i.e., reduced precision) of the memory by two distinct sources of noise. The first is described by a static, additive term (σ_1) that encompasses all potential one-time noise sources within a trial including noise associated with the sensory encoding and the motor response. The second is the dynamic degradation of memory precision over time that is modeled as the diffusion of the particle (*Figure 2a*). This diffusion corresponds to an increase in variability over time that is linear, with a slope equal to the diffusion constant (σ_2 ; *Figure 2b*). Consistent with past modeling studies (*Bays et al., 2009; Brady and Alvarez, 2015; Koyluoglu et al., 2017; Wei et al., 2012*), we accounted for the decrease in working-memory fidelity with item load by incorporating item noise over w, bnd FF00A1_4 1 ce ints-0.g ihe d5j0.1e 1 T1886 -1nd s

of averaging, and AtD produces a lower MN^2 and less variable responses than DtA. A summary of all framework variables can be found in *Table 1*.

To summarize, our two models describe two different possible ways for decision-relevant information to be stored in working memory prior to executing a decision. The different storage strategies result in different patterns of memory degradation, corresponding to trial-to-trial variability (imprecision) of decision reports that increase as a function of the length of the within-trial delay period. For

$p=0.029$). We also found these results were robust to uncertainty associated with model identifiability (participant-wise identifiability is given in *Figure 2—figure supplement 1*). Specifically, given different possible distributions of underlying strategy prevalence (proportions), the probability of obtaining the empirically observed distributions of models shown in *Figure 6a* for either set size while considering the average model identifiability was peaked near the observed strategy proportions. This result demonstrates that the observed proportions were not likely obtained due to misidentification-related chance. These probability distributions were also highly non-overlapping, which is consistent with a different prevalence of strategy use at the two different set sizes (*Figure 6b*).

These differences in strategy use did not correlate with the ages of the participants (Pearson correlation, *Figure 6—figure supplement 1*, $p>0.20$). These findings suggest that working-memory load might have affected our participants' decision strategies, such that a higher load corresponded to an increased tendency to discard information about individual samples (disk locations) and hold only the relevant computed decision variable in memory.

Sequential condition behavior

For the Sequential condition, we separately analyzed errors for Perceived reports of disks presented at the beginning (Early) or middle (Late) of a trial. Early Perceived reports tended to be relatively unbiased (two-sided t -test for H_0 : mean error=0, $p>0.05$; *Figure 7a*, full distributions in *Figure 7—figure supplement 1*; individual participant mean errors in *Figure 7—figure supplement 2a-d*) but became more variable over time in a roughly linear manner (*Figure 7d*), consistent with the predictions of the particle-

eLi

multiple quantities stored at once. Third, our DtA model also assumed that each item was stored individually. Alternatively, items could have been discarded or merged (chunked) (Krishnan *et al.*, 2018; Wei *et al.*, 2012), leading to different memory loads which could also affect performance. Fourth, most of our participants used strategies that were well described by the AtD or DtA model. However, under certain conditions (i.e., Sequential, set size 5) some participants seemed to use hybrid strategies. This kind of strategy would suggest extensive flexibility in when and how evidence is incorporated into computed decision variables, thereby placing potentially complex demands on working memory.

Both of our primary models were based on assumptions of a drifting memory representation. This random drift is traditionally associated with attractor models of working memory (Bays, 2014; Compte *et al.*, 2000; Macoveanu *et al.*, 2007; Wei *et al.*, 2012) that have been used extensively to describe the underlying neural mechanisms (Funahashi *et al.*, 1989; Shafi *et al.*, 2007; Takeda and Funahashi, 2002; Wimmer *et al.*, 2014). In these models, neural network activity is induced by an external stimulus and then maintained via excitatory connections of similarly tuned neurons and long-ranged inhibition. Random noise causes the center of this activity (which represents the stimulus) to drift in a manner that, dependent on the implementation, can depend on the delay duration, set size, and/or their interaction (Almeida *et al.*, 2015; Bays, 2014; Koyluoglu *et al.*, 2017). A recent implementation even can naturally compute a running average based on sequentially presented information (Esnaola-Acebes *et al.*, 2021). Our results imply that such models should be extended to support the flexible use of different strategies that govern when and how incoming information is used to form such averages. It will be interesting to see if such a flexible model can account for neural activity in the dorsolateral prefrontal cortex, which includes neurons with persistent activity that has been associated with both spatial working memory (Compte *et al.*, 2000; Constantinidis *et al.*, 2018; Riley and Constantinidis, 2016; Wei *et al.*, 2012; Wimmer *et al.*, 2014) and the formation of decisions based on an accumulation of evidence (Curtis and D'Esposito, 2003; Heekeren *et al.*, 2006; Heekeren *et al.*, 2008; Kim and Shadlen, 1999; Lin *et al.*, 2020; Philiastides *et al.*, 2011).

In conclusion, we found that in this spatial, continuous task, participant accuracy for both perceived and computed values was subject to working-memory limitations of both time and capacity. Additionally, we found behavior that was consistent with both the storage strategies we investigated. The fact that different participants employed different strategies for storing a computed value (such as a decision variable) and that these strategies have different consequences on overall accuracy has important implications for not only future neural network models of working memory, but also for future computational models of decision-making.

Materials and methods

Human psychophysics behavioral task

We tested 17 participants (4 males, 12 females, 1 chose not to answer; age range=22–87 years). The task was created with PsychoPy3 (Peirce *et al.*, 2019) and distributed to participants via [Pavlovia.org](https://pavlovia.org), which allowed participants to perform the task on their home computers after providing informed consent. These protocols were reviewed by the University of Pennsylvania Institutional Review Board (IRB) and determined to meet eligibility criteria for IRB review exemption authorized by 45 CFR 46.104, category 2.

Participants were instructed to sit one arm-length away from their computer screens during the experiment and to use the mouse to indicate choices. Each participant completed 1–2 sets of four blocks of trials in their own time.

The basic trial structure is illustrated in **Figure 1**. Each trial began with the presentation of a central white fixation cross (1% of the screen height). The participant was instructed to maintain fixation on this cross when not actively responding. The participant began each trial by placing the mouse over the cross and clicking, to allow for self-pacing and pseudo-fixation. Initiating a trial caused a white annulus of radius 25% of the screen height to appear. A block-specific memory array appeared 250 ms later, centered at an angle chosen uniformly and at random on the annulus. The array consisted of 1, 2, or 5 colored disks sized 1.5% screen in diameter. The angular difference between any two adjacent disks was at least 6°, and between the two most distal disks was at most 60°. The disks from clockwise to counter-clockwise were always presented in the same order: green, red, blue, magenta,

the effects of set size, delay duration, and task context on response variability using a two-way repeated measures ANOVA. On Simultaneous Perceived and Computed blocks, we used a 3 (delay duration: 0, 1, or 6 s) \times 3 (set size: 1, 2, or 5 disks) within-participant design. On Sequential Perceived blocks, we used a 2 (delay duration: 1 or 6 s) \times 3 (set size: 1, 2, or 5 disks) within-participants design for stimuli presented at the beginning of the trial (Early) and a 2 (delay: 0.5 or 3 s) \times 2 (set size: 2 or 5 disks) design for stimuli presented halfway through the trial (Late). On Sequential Computed blocks, we used a 2 (delay duration: 1 or 6 s) \times 3 (set size: 1, 2, or 5 disks) within-participants design. When the comparison included set size=1, data were always taken from the Simultaneous Perceived block.

To assess performance differences based on strategy use, additional analyses were performed once the data had been fit to the models and the best fit model had been selected (see below). These analyses included an assessment of response error variability in the Computed blocks using a 2 (model: AtD or DtA) \times 3 or 2 (delay: 0, 2, or 6 s Simultaneous condition, 1 or 6 s for Sequential) ANOVA with multiple comparisons to identify differences. To interrogate best fit parameter differences, two-sided *t*-tests were used to see if the mean difference in best-fit parameter between AtD and DtA participants was significantly different from 0 for both Simultaneous and Sequential conditions. To assess learning effects, a two-sided, paired *t*-test was used to see if the mean or standard deviation of error responses in set size 5 Sequential conditions differed between the first and second half of trials (we found no difference at either delay: for 1 s delay $p=0.67$ and 0.11 for mean and standard deviation, respectively; for 6 s delay $p=0.75$ and 0.98 for mean and standard deviation, respectively).

Model-based analyses

Our models were based on principles of working memory that are well described by bump-attractor network models (Compte et al., 2000; Laing and Chow, 2001; Wimmer et al., 2014). In such models, stimulus location is represented by a ‘bump’ in activity from neurons tuned to that and similar locations. These neurons recurrently activate each other, maintaining a bump of activity even after stimulus cessation. However, because of the stochastic nature of neural activity and synaptic transmission (Faisal et al., 2008), there is variability in which neurons have the most activity at any given time (and thus are the center of the bump representing the stimulus). This variability in bump center corresponds to variability in the location representation and a degradation of the memory representation over time. The dynamics of this bump can be described as a diffusion process that obeys Brownian motion (Compte et al., 2000; Kilpatrick, 2018; Kilpatrick et al., 2013; Laing and Chow, 2001). We used this simplified description in our models as follows.

Perceived values in working memory

A single point (i.e., the central spatial location of a single disk), x_1 , is assumed to be represented in working memory by $x_{t,1}$, where t represents the time since the removal of the stimulus. We assume that $x_{t,1}$ evolves like a sample from a Brownian-motion process. Specifically, when x_1 is observed, it is encoded with some perceptual noise, σ_p . Therefore, at time zero, $x_{0,1} \sim N(x_1, \sigma_p^2)$. This representation accumulates noise over time with some diffusion constant, σ_d^2 , further degrading the representation of $x_{t,1}$ from x_1 such that $x_{t,1} \sim N(x_1, \sigma_p^2 + t \sigma_d^2)$. There is additional motor noise in the participant’s report, $r_{t,1}$, and we denote the variance of this motor noise by σ_m^2 . Mathematically, it is equivalent to add the motor noise at the beginning or the end of the diffusion of $x_{t,1}$ when considering the report, $r_{t,1}$. In our model, we thus represent the sum of the perceptual and motor noise as a single, static noise term. Hence, we show simulated trajectories of as(i80n200D>>> BDC 18.659 -00 9 267.2879 2i /MC8 ()TJEMC ET/P

e

corresponding to the increased memory load. The representation of the Late item then diffuses for only half of the delay time, T (see **Figure 2d and e**). We formalized this process with the following model for the report error of the Early ($e_{T,NE}$) and Late ($e_{T,NL}$) items:

$$e_{T,NE} \sim \mathcal{N}(0, \eta_{NE} + \pi^2 \sigma_1^2 (N-1)^A + \pi^2 \sigma_1^2 N^A)$$

following conditions. Perceived: delays 1, 3, and 6 s; array size 1 (*Equation 3a*). Perceived: delays 3 and 6 s, array size N for both Early (*Equation 6a*) and Late (*Equation 6b*) items. Computed: delays 3 and 6 s, array size N (*Equation 7* for AtD or *Equation 8* for DtA).

Because the mean error for each individual participant was not always 0, when fitting the AtD and DtA models we used the empirical mean error from the condition being fitted as a fixed bias term in the model. Mean error and CIs for each participant for each condition are shown in *Figure 3—figure supplements 2 and 3; Figure 3—figure supplements 2 and 3*.

We obtained separate maximum-likelihood fits for AtD and DtA models for each individual participant, using the function `fmincon` in MATLAB to minimize the summed negative log-likelihood of obtaining the observed errors for a given condition according to the above equations. Initial parameter values were randomized and the fitting repeated to avoid local minima. Because all models within a given condition had the same number of parameters, we compared log-likelihoods to determine the best-fitting model for a given participant. Because the number of parameters is the same, comparing likelihoods produces equivalent model selection to BIC or AIC.

Assessing model assumption and identifiability

To assess how well each participant's data matched the assumptions of the AtD and DtA models, we also fit a line to the variances of response errors across delays for a given condition for a given participant to obtain empirical estimates of the various diffusion constants (e.g., slope of lines in *Figure 2b*; empirical estimate of a Perceived value,

relationships. Participants whose empirical diffusion constant relationships fell within the central 95% of the simulated expected range were considered well fit by their model.

To assess model identifiability, for each participant and condition, we fit both models to the results of each set of 1000 simulations generated using the best-fitting parameters from the best-fitting model for that participant and condition. We used the log-likelihoods to determine the best model for each simulation and determined the percentage of correctly identified models. We used these models as each pd-

Wimmer K, Nykamp DQ, Constantinidis C, Compte A. 2014. Bump attractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nature Neuroscience*