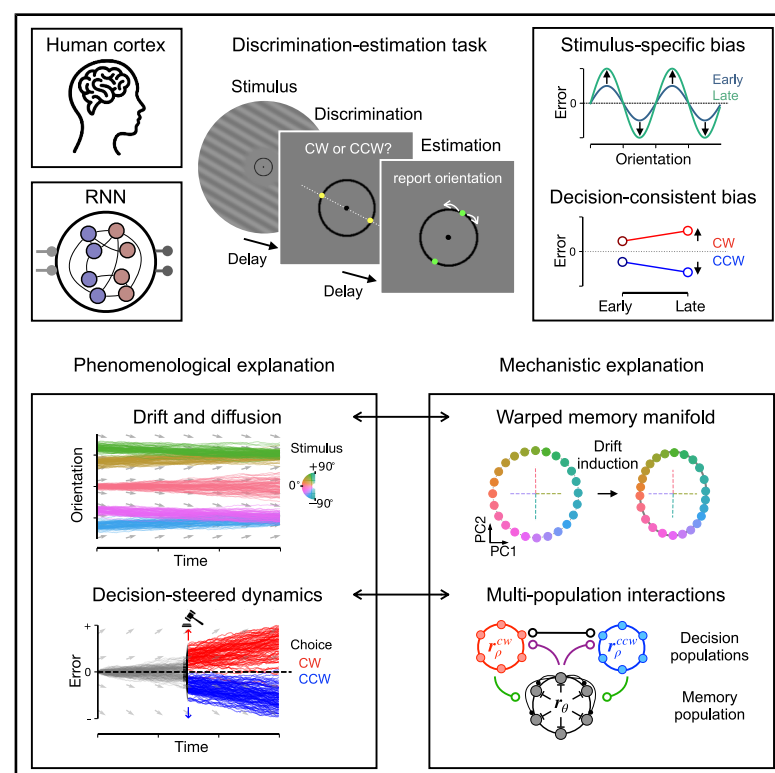


Attractor dynamics of working memory explain a concurrent evolution of stimulus-specific and decision-consistent biases in visual estimation

Graphical abstract



Authors

Hyunwoo Gu, Joonwon Lee, Sungje Kim, ..., Jun Hwan (Joshua) Ryu, Sukbin Lim, Sang-Hun Lee

Correspondence

sukbin.lim@nyu.edu (S.L.),
visionsl@snu.ac.kr (S.-H.L.)

In brief

People exhibit biases when perceiving features of the world, shaped by both external stimuli and prior decisions. By tracking behavioral, neural, and mechanistic markers of stimulus- and decision-related biases, Gu et al. show that working memory's attractor dynamics—instantaneously updated by decisions—comprehensively explain the growth and interplay of these biases.

Highlights

- Delayed estimates are biased in stimulus-specific and decision-consistent ways
- Behavioral and neural growth of stimulus-specific bias implies drift to attractors
- Decisions steer stimulus-specific memory drift, augmenting decision-consistent bias
- Task-optimized RNNs reveal mechanisms for such decision-steered attractor dynamics

Article

Attractor dynamics of working memory explain a concurrent evolution of stimulus-specific and decision-consistent biases in visual estimation

Hyunwoo Gu,^{1,2,3} Joonwon Lee,¹ Sungje Kim,¹ Jaeseob Lim,¹ Hyang-Jung Lee,¹ Heeseung Lee,^{1,7} Min Jin Choe,¹ Dong-gyu Yoo,¹ Jun Hwan (Joshua) Ryu,^{2,3} Sukbin Lim,^{4,5,6,*} and Sang-Hun Lee^{1,8,*}

¹Department of Brain and Cognitive Sciences, Seoul National University, 1 Gwanak-ro, Seoul 08826, Republic of Korea

²Department of Psychology, Stanford University, Stanford, CA 94305, USA

³Wu Tsai Neurosciences Institute, Stanford University, Stanford, CA 94305, USA

⁴Shanghai Frontiers Science Center of Artificial Intelligence and Deep Learning, NYU Shanghai, 567 West Yangsi Road, Shanghai 200126, P.R. China

⁵Neural Science, NYU Shanghai, 567 West Yangsi Road, Shanghai 200126, P.R. China

⁶NYU-ECNU Institute of Brain and Cognitive Science, NYU Shanghai, 3663 Zhongshan Road North, Shanghai 200062, P.R. China

⁷Department of Psychological and Brain Sciences, Dartmouth College, Hanover, NH 03755, USA

⁸Lead contact

*Correspondence: sukbin.lim@nyu.edu (S.L.), visionsl@snu.ac.kr (S.-H.L.)

<https://doi.org/10.1016/j.neuron.2025.07.003>

SUMMARY

Sensory evidence tends to be fleeting, often unavailable when we categorize or estimate world features. To overcome this, our brains sustain sensory information in working memory (WM). Although keeping that information accurate while acting on it is vital, humans display two canonical biases: estimates are biased toward a few stimuli (“stimulus-specific bias”) and prior decisions (“decision-consistent bias”). Integrative—especially neural mechanistic—accounts of these biases remain scarce. Here, we identify drift dynamics toward discrete attractors as a common source of both biases in orientation estimation, with decisions further steering memory states. Behavior and neuroimaging data reveal how these biases co-evolve through the decision-steered attractor dynamics. Task-optimized recurrent neural networks suggest neural mechanisms that enable categorical decisions to emerge from WM for continuous stimuli while updating their trajectory, warping decision-consistent biases under stimulus-specific drift.

INTRODUCTION

Adapting to our surroundings engages us in various perceptual tasks on the same object,^{1,2} like judging its categorical state (e.g., whether an apple is “small” or “large”) and estimating its exact state (e.g., “precise size” of an apple). These tasks often occur in succession, with cognitive processes for earlier tasks influencing later ones. “Decision-consistent bias” is a prime example, where our estimate of a feature aligns with the categorical state of our previous decision (e.g., after deciding on “large,” size estimates tend to be larger than the actual size). Understanding this can provide insights into the brain’s flexible use of feature representations under varying task demands. While research has clarified how categorical decisions are formed from sensory evidence,^{3–7} how the decision-forming process influences the subsequent retention of sensory evidence for future reuse remains unclear.

Decision-consistent bias, once viewed as a perceptual illusion caused by biased readouts of unbiased sensory representations,² has been recently reconceptualized^{8–11} as involving

post-perceptual processes. In these studies,^{2,8–11} since sensory inputs from target stimuli are no longer available during a subsequent estimation task, sensory evidence must be held in working memory (WM).¹² However, efforts to explain decision-consistent bias in the context of WM have been surprisingly scarce, especially given WM’s dynamic nature^{13–15} and its close relationship with decision-making (DM).^{15–18}

Alongside decision-consistent bias, perceptual estimation exhibits another prominent bias: estimates tend to cluster around specific points in a feature space (e.g., position estimates deviate from cardinal to oblique meridians). This phenomenon, called “stimulus-specific bias,” is common across various domains, including spatial position,^{19,20} motion direction,^{21,22} orientation,^{23,24} and color.^{25,26} Although Bayesian approaches have yielded normative accounts of why stimulus-specific bias occurs,^{27–29} our understanding of its neural-mechanistic origins remains limited. Moreover, its relationship with decision-consistent bias has not been explored, despite both occurring in similar estimation tasks. These gaps underscore the need for an integrated, mechanistic-level account to clarify their coexistence

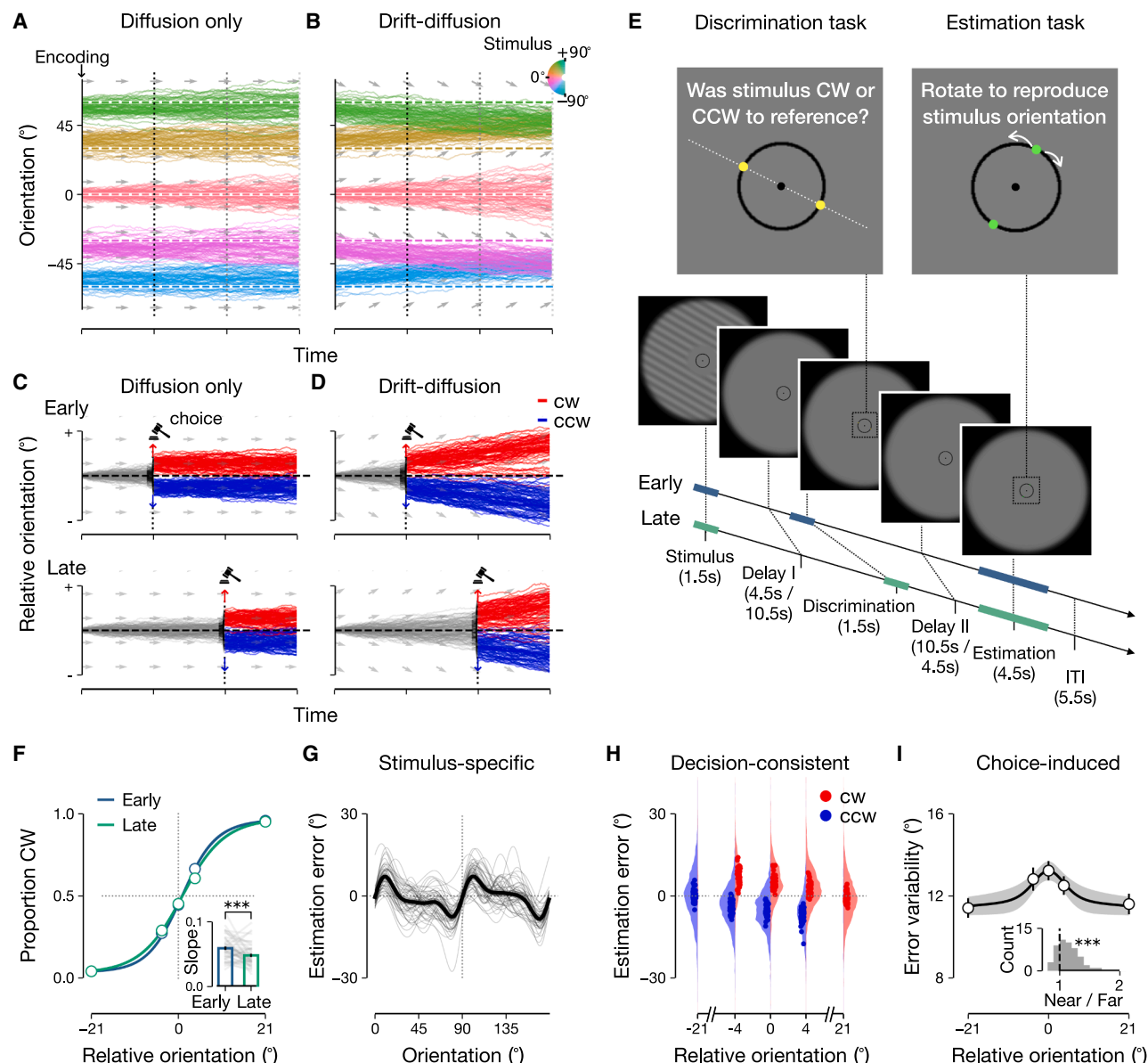


Figure 1. Probing the impact of WM dynamics on behavioral biases

(A–D) Single-trial memory trajectories simulated under diffusion-only (A and C) and drift-diffusion (B and D) dynamics. Gray arrows indicate drift directions. (A and B) Impact on stimulus-specific bias. (C and D) Impact on decision-consistent bias. Decision timing (gravel) varies between top (early) and bottom (late) panels. (E) Task paradigm. Top: yellow and green dots serve as the reference for discrimination and the response frame for estimation, and dotted lines, arrows, and white text are shown for illustration purposes. Bottom: blue and teal bars represent the task epochs in early- and late-DM trials.

(F) Psychometric curves of discrimination, pooled across individuals. Late-DM shows shallower slopes (inset: thin lines, individuals; error bars, SEM; paired t test, $p = 0.0008$).

(G–I) Behavioral signatures of stimulus-specific (G), decision-consistent (H), and choice-induced (I) biases. (G) Estimation errors captured via von Mises function derivatives: thin, individuals; thick, mean across participants \pm SEM. (H) Choice-conditioned distributions of estimation errors: dots, means for individuals; patches, pooled densities. (I) Error variability measured as interquartile range: circles with error bars, mean \pm SEM; curve with a gray shadow, Gaussian fit to the data with SEM; inset, variability ratio histogram, with a significant tendency above 1, Wilcoxon signed-rank test, $p < 10^{-5}$. *** $p < 0.001$.

and interaction in a dynamic context. Focusing on only one bias might ignore the interconnected roots linked to the other.

We hypothesized two intrinsic dynamics of WM that critically influence the co-evolution and interaction of the two biases. Memory states may gradually and randomly shift in a feature

space, engendering diffusive representations that grow noisier yet remain unbiased^{13,14,30} (Figure 1A). Alternatively, memory states may not only diffuse but also drift toward a few stable points (attractors), engendering drift, and diffusive representations that become increasingly biased and noisy^{26,31–33}

(Figure 1B). Assuming decision-formation processes steer ongoing memory states in the choice-consistent direction (Figures 1C and 1D, arrows), we expect that the two biases will undergo different time courses depending on whether WM dynamics are driven by diffusion alone (Figure 1C) or by both diffusion and drift (Figure 1D), and when decisions are made (Figures 1C and 1D, top versus bottom panels).

To examine these expectations, we designed a task in which participants made sequential categorical choices and point estimates about a remembered stimulus orientation, with varying intervening delays (Figure 1E). Monitoring human individuals' task performance through their choices and estimates while decoding orientation memory from functional magnetic resonance imaging (fMRI) of their visual cortex, we tracked the behavioral and neural signatures of both biases over a prolonged delay period. To further examine whether simple neural mechanisms could recapitulate WM and DM interaction in humans, we also trained recurrent neural networks (RNNs) on the same task and analyzed their dynamics.

Our results across behavioral, fMRI, and RNN analyses convergently indicate “decision-steered attractor dynamics of WM” as the core mechanism underlying the co-evolution of stimulus-specific and decision-consistent biases. Memory states drift toward attractors, influencing categorical decisions, which in turn bias memory trajectories, creating cascading effects that amplify both biases. RNN simulations mirrored human behavior and fMRI dynamics, showing that decisions emerge through modular interactions among three populations: a WM population maintaining orientation memory and two DM populations competing for choices, with feedback from DM populations steering WM attractor dynamics. Our work provides an integrated neural-mechanistic explanation of the fundamental biases in perceptual estimation and their evolving interaction, which have not been previously addressed.

RESULTS

Behavioral signatures of stimulus-specific and decision-consistent biases

During the fMRI scan, participants memorized the orientation of a briefly shown grating and reported it after a 16.5-s delay. To probe the impact of DM on WM, they first performed a discrimination task during the delay (Figure 1E, top). In this task, a dot pair around the fixation (“reference”) appeared, and participants decided, under moderate time pressure (1.5 s), whether the remembered orientation was tilted clockwise (CW) or counter-clockwise (CCW) relative to the reference, whose angle was randomly determined. The timing of the discrimination task varied, occurring either 4.5 or 10.5 s after stimulus offset (Figure 1E, bottom). As anticipated from the temporal deterioration of WM,^{15,18} discrimination performance declined when tested later, as indicated by a shallower psychometric curve (Figure 1F). In the estimation task, participants rotated another dot pair (“report frame”) to match the remembered orientation, starting from a randomly chosen angle within 180°.

We confirmed stimulus-specific bias in the estimation task: estimates were repelled from cardinal orientations and attracted toward oblique orientations (Figure 1G). This well-known phe-

nomenon, called “cardinal repulsion,”^{23,27,34,35} showed modest variation in shape and size among individuals.

Decision-consistent bias was also evident when estimation errors were conditioned on discrimination choices, deviating from the reference in line with the choice. This phenomenon, also known as “reference repulsion,”^{2,9} was pronounced in trials where the reference orientation was close to the stimulus (-4° , 0° , and 4° on the x axis in Figure 1H). As noted previously,⁹ if choices induce a bias in estimates, the marginal error distribution must widen as the stimulus and reference become more similar in orientation (see Figures S1D–S1K for the rationale). Consistent with this, error variability was greater in near-reference (relative orientation $\hat{\theta} \in \{-4^\circ, 0^\circ, 4^\circ\}$) trials compared with far-reference ($\hat{\theta} \in \{-21^\circ, 21^\circ\}$) trials (Figure 1I). This increased variability was confirmed across various datasets when the reference was relevant to decisions but not when it served as a distractor (Figures S1A–S1C). We will call this bias “choice-induced bias,” to distinguish it from decision-consistent bias. The latter refers to any deviation aligned with a choice,^{2,9,10} measurable from any joint observations of discrete choices and continuous estimates (Figure 1H). Choice-induced bias is a particular kind of decision-consistent bias, where commitment to a categorical choice influences estimates beyond what statistical conditioning would predict⁸ (Figure 1I).

Predicting how WM dynamics affect the time courses of biases

We developed phenomenological models to predict how WM dynamics influence bias time courses, with minimal assumptions about sensory encoding and decision commitment's impact on memory. One model incorporates only diffusion dynamics, allowing random shifts of memory states across trials without bias (Figures 2A–2D), while the other includes additional drift dynamics that systematically drive memory states in specific directions (Figures 2E–2H).

Both models assumed initial orientation memory states follow the efficient encoding principle,^{27–29} which allocates more resources to frequently encountered stimuli,³⁶ leading to higher encoding precision around cardinal orientations. This results in specific error patterns: a mean repulsion away from cardinal orientations and increased variance around oblique ones, aligning with behavioral data. These biases are embedded in the sensory input to the WM system, so memory states already show stimulus-specific bias from the start (Figures 1A, 1B, 2B, and 2F, darkest curves). To incorporate choice-induced bias into dynamic WM, we also assumed that memory states are abruptly shifted in the chosen direction by a constant amount during the discrimination epoch (vertical arrows in Figures 1C and 1D). This pulse-like shift expands the marginal error distribution for near-reference trials, pushing choice-conditioned distributions apart (Figures 2A and 2E).

Based on these assumptions, we predicted how WM dynamics influence biases in both the diffusion-only and drift-diffusion models. If diffusion solely governs the dynamics (Figures 1A and 2A), stimulus-specific bias remains unchanged (Figure 2B) because stochastic fluctuations of memory states with a zero mean do not cause any systematic deviations (gray arrows in Figure 1A). Conversely, decision-consistent bias varies with

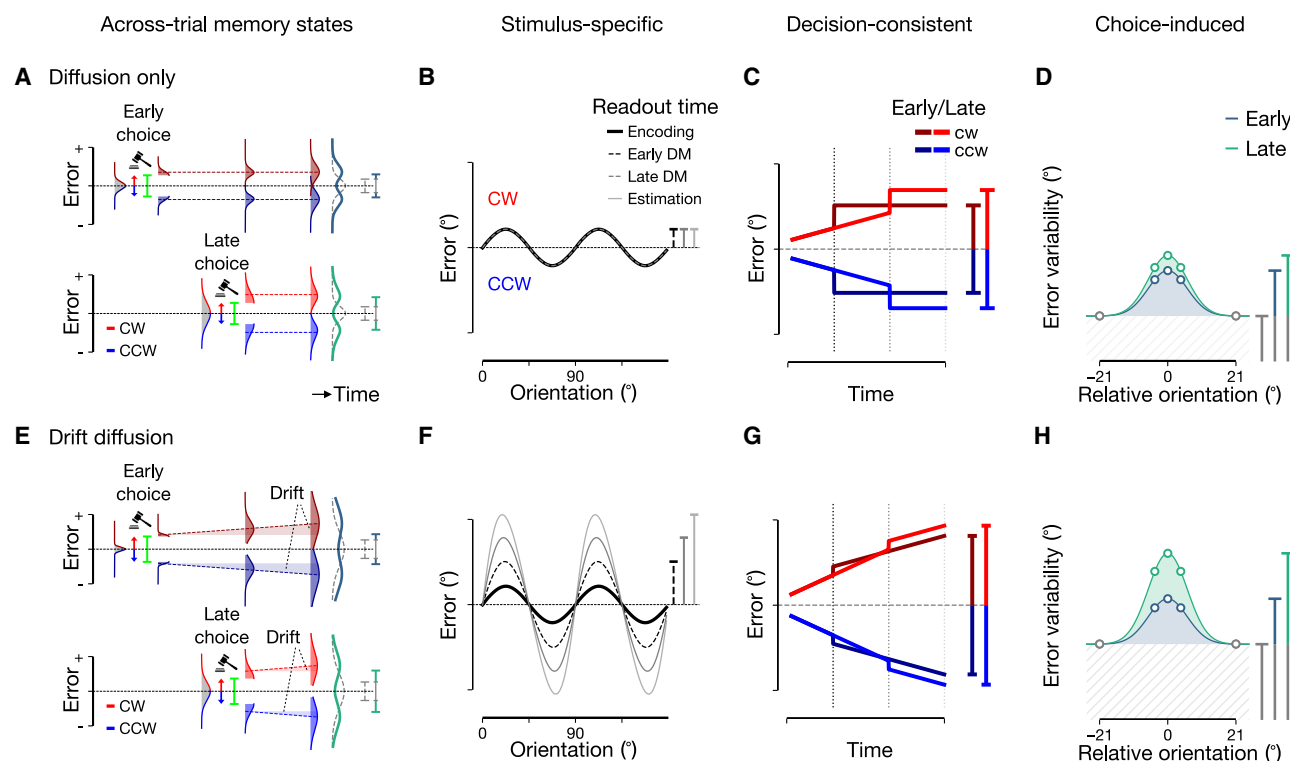


Figure 2. Model predictions of biases under WM dynamics

(A–D) Diffusion-only model.

(E–H) Drift-diffusion model. (A and E) Choice-conditioned memory distributions for early (top) and late (bottom) DM trials, with conditional means indicated by horizontal dashed lines. Colored arrows and lime markers represent choice-induced bias during discrimination. Green densities represent marginal error distributions during estimation, alongside dashed densities representing those without choice-induced bias. (B and F) Stimulus-specific biases across different task epochs. (C and G) Decision-consistent biases for early- and late-DM trials. (D and H) Choice-induced bias captured by near-reference variability for early- and late-DM trials, alongside dashed gray regions representing error variability without choice-induced bias.

decision timing: as previous studies^{8,9} indicate, stochastic sensory fluctuations, when conditioned on a choice, contribute to this bias. When diffusion governs WM dynamics, the distribution of memory states broadens over time. Thus, more delayed decision timing leads to greater separation of memory states associated with different choices (colored dashed lines in Figure 2A), resulting in an increased decision-consistent bias (Figure 2C).

Next, suppose both diffusion and drift govern WM dynamics (Figures 1B and 2E). Then, both biases vary with decision timing. A recent study on WM for color²⁶ suggests that drift causes memory states to approach discrete attractors, leading to increased stimulus-specific bias over time, while many others^{27–29,37–41} attribute this bias to sensory encoding. Similar drift dynamics may also govern WM for orientation (gray arrows in Figure 1B), causing stimulus-specific bias to grow over time (gray vertical bars in Figure 2F). As previously noted, decision-consistent bias will also increase with decision timing due to WM diffusion (colored vertical bars in Figure 2G).

Furthermore, diffusion dynamics predict that the broadening of the estimate distribution in near-reference trials—indicating choice-induced bias—will be more pronounced in the late-DM condition than in the early one (Figures 2D and 2H). This is because, despite equal choice-induced bias in both conditions

(colored arrows in Figures 2A and 2E), its effect on distribution expansion intensifies with decision delay under diffusion dynamics (marginal distributions in Figures 2A, 2E, and S1K).

In summary, both diffusion-only and drift-diffusion models predict that, with decision delay, decision-consistent bias and estimation variability increase in near-reference trials (Figures 2C, 2D, 2G, and 2H). However, they differ in stimulus-specific bias: it stays constant over time in the diffusion-only model but increases in the drift-diffusion model (Figures 2B and 2F).

Growth and decision-timing dependency of behavioral biases

To determine whether stimulus-specific bias increases during the delay, we compared the bias magnitude at the early (4.5 s post-stimulus) and late (10.5 s post-stimulus) discrimination epochs. Since direct estimation errors were unavailable during discrimination, we inferred the bias from the psychometric curve for each stimulus orientation (Figure 1F), using the deviation of subjective equality from the actual orientation as a proxy. For this, we fitted the bias weight parameter, assuming the bias varies in magnitude but retains its shape (see STAR Methods). We found that the bias increased over time, with a significantly

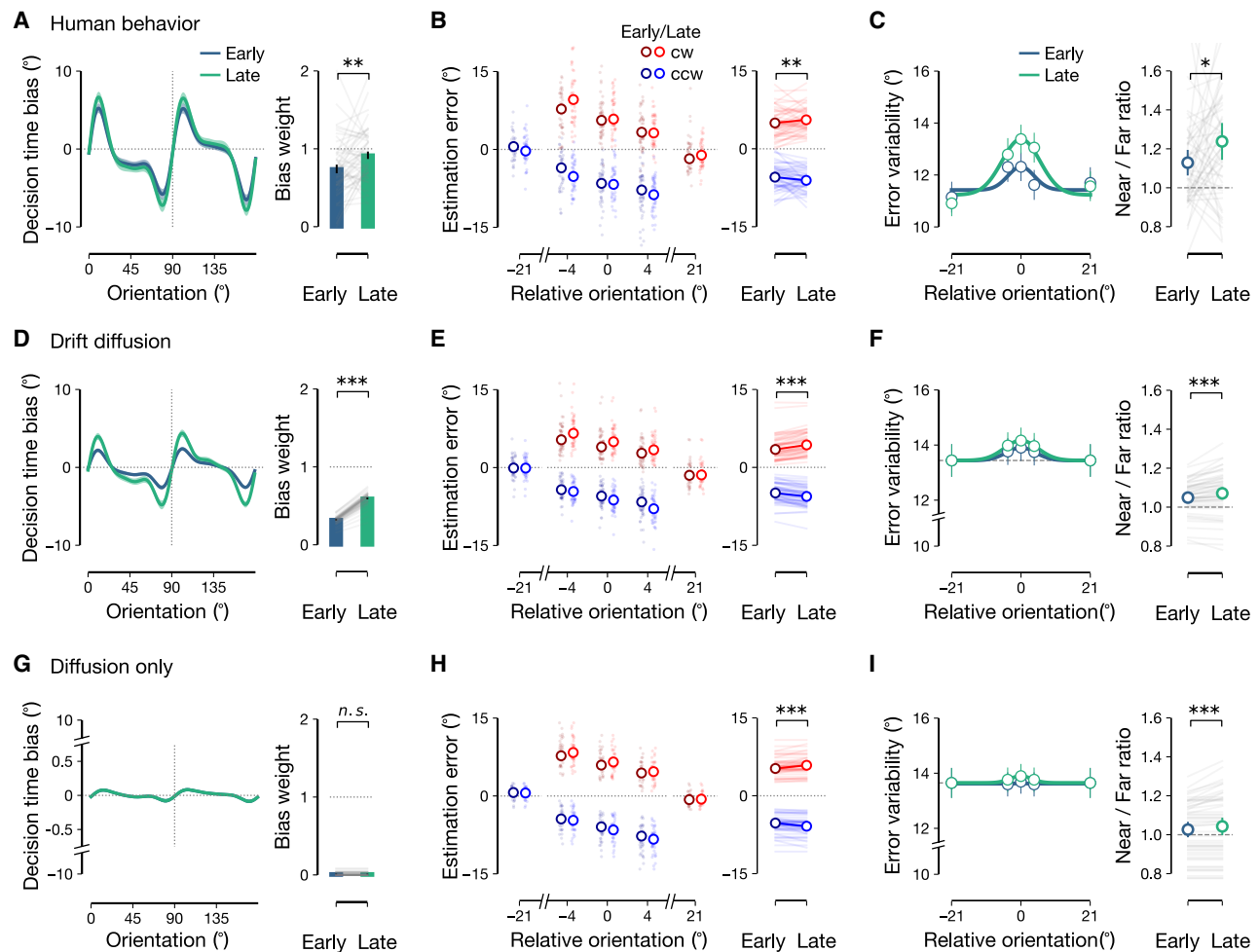


Figure 3. Biases in human behavior and models

(A–C) Human data. (A) Stimulus-specific bias at discrimination: left: bias estimates across orientations (lines with shades, means \pm SEMs across individuals); right: bias weights relative to estimation bias (bars, mean; lines, individuals; error bars, \pm SEM; paired t test, $p = 0.0032$). (B) Decision-consistent bias in estimation (left) and their averages across near-reference trials (right): dots and thin lines, individuals; circles, medians \pm SEM; paired t test, $p = 0.0052$. (C) Error variability for early- and late-DM (left) and near-/far-reference variability ratio (right): circles, across-individual means; lines, individuals; Wilcoxon signed-rank test, $p = 0.0450$. (D–I) Biases simulated by drift-diffusion (D–F) and diffusion-only (G–I) models. Format as in (A–C). Paired t test, $p < 10^{-10}$ (D), $p < 10^{-10}$ (E), $p = 0.3614$ (G), $p = 0.0001$ (H); Wilcoxon signed-rank test, $p < 10^{-5}$ (F), $p = 0.0007$ (I). *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$; ns $p > 0.05$.

greater bias at the late compared with the early decision (Figure 3A).

In contrast to their distinct predictions for stimulus-specific bias, both models predicted similar patterns of decision-consistent bias and near-reference variability (Figures 2C, 2D, 2G, and 2H). Consistent with these, both measures were greater in the late-DM condition (Figures 3B and 3C). To test the prediction regarding decision-consistent bias, we compared the reported orientation between CW-choice and CCW-choice trials for the near-reference condition, and the difference was significantly larger in the late than the early DM condition (Figure 3B). For near-reference variability, we examined the difference in error variance between choice-unconditioned near-reference and far-reference trials, and the difference was also greater in the late-DM condition (Figure 3C).

To quantify drift dynamics, we fitted the models—with and without drift—to behavioral data (see STAR Methods). The model with drift outperformed the one without drift, according to the Bayesian information criterion (BIC) in conjunction with cross-validated log likelihoods (Figure S2A). Drift rate parameters indicate memory states drift modestly (less than $1^\circ/\text{s}$ in median; w_K in Figure S2C). In *ex-post* simulations, the drift model accurately reproduced the observed growth and shape of stimulus-specific bias (Figures 3D, S3A, and S3B), unlike the non-drift model (Figures 3G and S3C). Both models captured the observed decision-consistent bias patterns and variability increases in near-reference trials (Figures 3E, 3F, 3H, and 3I), but a model solely based on efficient coding, without drift or diffusion dynamics, failed to do so (Figures S3D–S3F).

In summary, incorporating drift-and-diffusion dynamics with transient shifts in the chosen direction into WM effectively explains the growth of stimulus-specific and decision-consistent biases observed in our task.

Influence of stimulus-specific drift on decision-consistent bias before and after DM

Our analyses indicate that memory states consistently drift toward attractors, with decision-consistent bias increasing with decision timing. This implies that decision-consistent bias follows distinct time courses depending on the stimulus's position relative to these attractors in orientation space. We will first formalize this implication using our earlier drift-diffusion model, then test it against human data.

This implication involves dividing decision-consistent bias into its pre-decision and post-decision components. The pre-decision bias (b_{pre}) refers to the difference in mean between the choice-conditioned distributions of memory states present at the onset of discrimination, whereas the post-decision bias (b_{post}) develops after discrimination until estimation (see [STAR Methods](#) and [Method S1.1](#) for definitions).

Under diffusion-only dynamics, b_{pre} is expected to be greater in the late-DM condition than the early DM condition, due to increased separation between choice-conditioned distributions ($\Delta b_{pre} > 0$; dark gray patch in [Figure 4A](#)). Conversely, b_{post} should remain constant in both conditions if choice-induced bias remains unchanged in magnitude, as assumed ($\Delta b_{post} = 0$; light gray patch in [Figure 4A](#)). These predictions were confirmed by the b_{pre} and b_{post} values derived from the *ex-post* simulation of the diffusion-only model ([Figure 4D](#)).

However, under drift-diffusion dynamics, the two components of decision-consistent bias display distinct, stimulus-dependent patterns. Consider the case where drift, diffusion, and choice-induced bias are moderate, so that memory states starting far from the attractors do not reach them during the delay, consistent with the best-fit parameters from the behavioral data ([Figure S2](#)). For orientations positioned between the attractors (e.g., cardinal orientations), memory states diverge from the stimulus (pink lines around 0° in [Figure 1B](#)), increasing the separation of choice-conditioned memory distributions over time, driven by the congruence between the diverging direction and the decision-consistent direction ([Figure 4B](#)). This leads to a greater b_{pre} in the late versus early DM condition ($\Delta b_{pre} > 0$; dark gray patch in [Figure 4B](#)), while b_{post} should be smaller in the late-DM condition ($\Delta b_{post} < 0$; light gray patch in [Figure 4B](#)). Conversely, near the attractors (e.g., oblique orientations), converging drift decreases the separation of choice-conditioned memory distributions over time, counteracting the choice-conditioned separation ([Figure 4C](#)).

The drift-diffusion model also implies a covariation across individuals. While the best-fit drift rates are moderate (maximum $w_K < 2.5^\circ/\text{s}$ in [Figure S2C](#)), their differences across individuals predict systematic changes in decision-timing-dependent biases. Specifically, higher drift rates lead to more pronounced changes in both Δb_{pre} and Δb_{post} . Thus, given their opposite signs and dependence on drift rate, a negative correlation between Δb_{pre} and Δb_{post} is predicted across individuals, driven by drift rate variability ([Figures S3G and S3H](#)).

These two implications were confirmed by the *ex-post* simulation data from the drift-diffusion model ([Figure 4E](#)), along with human data analyses that did not rely on model estimates ([Figure 4F](#)). First, the Δb_{pre} and Δb_{post} were positive and negative, respectively, for diverging drift orientations (pink in right panels), but both were near zero for converging drift orientations (green in right panels). The joint distributions of Δb_{pre} and Δb_{post} were significantly separated. Here, considering the individual differences in stimulus-specific bias shape ([Figure 1G](#)), we determined whether a converging or diverging drift governs a given orientation based on participant-specific bias curves (see [STAR Methods](#)). Second, when the data were pooled across orientations ([Figures 4E and 4F](#), left panels), the Δb_{pre} and Δb_{post} were still positive and negative, respectively. This was anticipated because the Δb_{pre} and Δb_{post} were large near diverging drift orientations. Across individuals, the Δb_{pre} and Δb_{post} were negatively correlated ([Figure 4F](#)).

In summary, phenomenological models reveal how the drift-diffusion dynamics of WM intricately shape decision-consistent bias before and after decisions in a stimulus-specific way, supported by human behavior.

Drift-diffusion dynamics in cortical signals of orientation memory

The behavioral analysis examined the biases at snapshot moments of discrimination and estimation. To verify and expand on these findings beyond these moments, we decoded the WM signal of stimulus orientation from the blood-oxygenation-level-dependent (BOLD) measurements via inverted encoding analysis,^{42–44} tracking the biases over time in that decoded signal. Focus was primarily on early visual areas, V1, V2, and V3, given their high-fidelity WM for orientation,^{43,45,46} with parietal and frontal areas included for comparison ([Figures S4E–S4I](#)).

As implied by drift-diffusion dynamics, the cortical signal of orientation memory confirmed the growth of stimulus-specific bias over time. Stimulus orientation in WM was decodable from the early visual cortex with significant fidelity across all trial time points ([Figure S4D](#)), unattributable to eye movement confounds ([Figure S5](#)). Its trajectories, conditioned on stimulus orientations, drifted away from cardinal and toward oblique orientations ([Figure 5A](#)). To quantify these attractor dynamics, we tracked bias strength using linear regression weights that related each time point's bias to each individual's behavioral stimulus-specific bias (thin gray curves in [Figure 1G](#)). The bias weight was initially low, consistent with the efficient coding framework^{28,29,37} ([Figure S7](#)), and increased to match those observed in the behavioral errors ([Figure 5B](#)). Additionally, representational similarity analyses⁴⁷ and simulated population responses with heterogeneous tuning curves ([Figure S6](#); [Method S2.1](#)) confirmed the growth of stimulus-specific bias in memory representations.

We note two caveats in inferring WM dynamics from BOLD signals. First, BOLD signals may appear to change more gradually than actual neural activity due to hemodynamic effects. Second, the limited temporal resolution of BOLD signals can cause interference between stimulus and reference orientations during the discrimination epoch. Indeed, the brain signals transiently shifted toward the near-reference orientation, especially in the

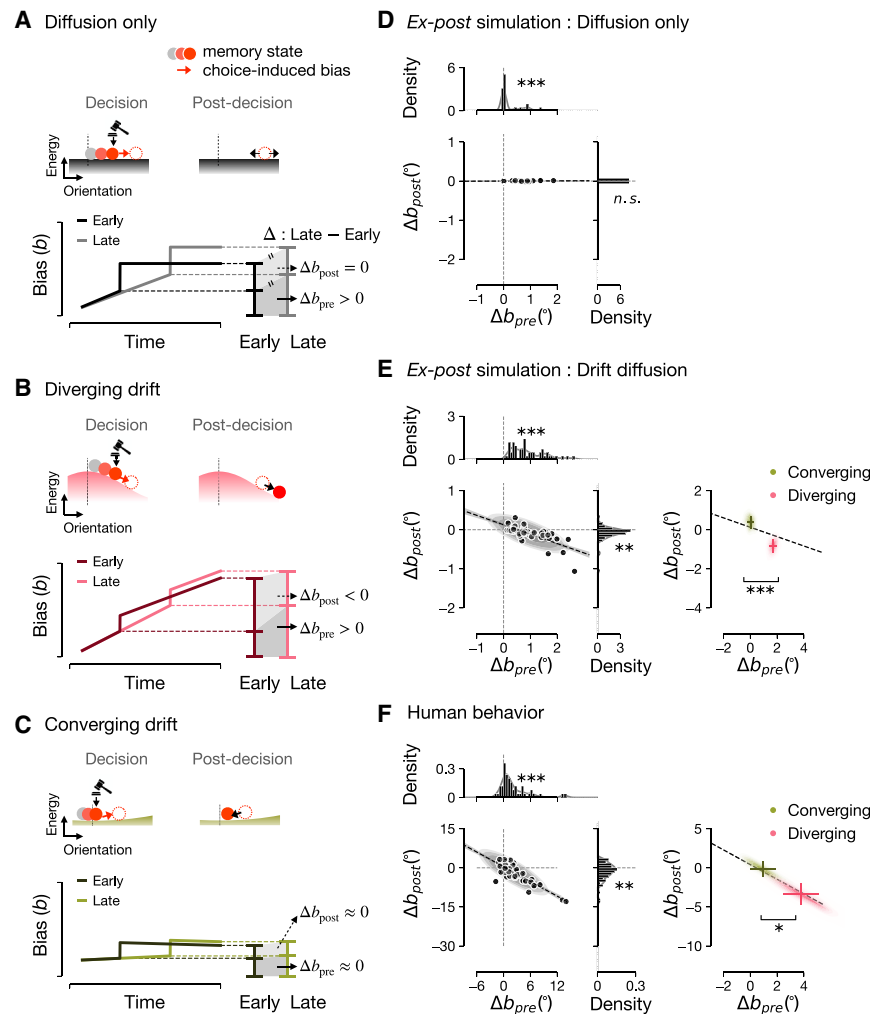


Figure 4. Decision-consistent biases before and after DM

(A–C) Schematics of pre- and post-decision biases under diffusion-only (A), diverging-drift (B), and converging-drift (C) WM dynamics. Top: memory dynamics on energy landscapes: gray to dark red circles, early to late states; dotted circles, states shifted by DM; red arrows, choice-induced biases; black arrows, post-decision biases. Bottom: decision-consistent bias trajectories: dark/light lines, early/late-DM conditions.

(D) Across-individual Δb_{pre} and Δb_{post} values simulated by the diffusion-only model, shown as joint (left bottom) and marginal (top, right) distributions: dashed lines with shades, regression lines; dots, individuals.

(E) Simulations by the drift-diffusion model. Format as in (D), except for the right panel, where across-individual averages of Δb_{pre} and Δb_{post} are shown separately for diverging and converging orientations.

(F) Human data, with format as in (E).

Signs of Δb_{pre} and Δb_{post} were tested with one-sample t test (Δb_{pre} , $p < 10^{-4}$, Δb_{post} , $p = 0.5924$ in D; Δb_{pre} , $p < 10^{-10}$, Δb_{post} , $p = 0.0024$ in E; Δb_{pre} , $p < 10^{-4}$, Δb_{post} , $p = 0.0022$ in F). Correlations were Pearson's coefficients ($r = 0.203$, $p = 0.1568$ in D; $r = -0.726$, $p < 10^{-8}$ in E; $r = -0.829$, $p < 10^{-10}$ in F). Distances along regression lines were measured (1.941° , $p < 10^{-4}$ in (E); 4.215° , $p = 0.0457$ in F). *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$; ns $p > 0.05$.

(Δb_{pre} and Δb_{post}) matched the predictions from the drift-diffusion model (Figure 5E), showing positive Δb_{pre} , negative Δb_{post} , negative correlation between Δb_{pre} and Δb_{post} , and distribution

separation between Δb_{pre} and Δb_{post} for the stimulus orientations with converging and diverging drifts.

Lastly, to assess how accurately cortical signals reflect behavioral biases, we convolved the predicted memory trajectories from the best-fit diffusion-only and drift-diffusion models with the canonical HRF (see STAR Methods and Figures S4A and S4B). The stimulus-specific trajectories predicted by the drift-diffusion model (Figures S7A and S7B) closely matched the BOLD signals of orientation memory (Figures 5A and 5B), unlike those from the diffusion-only model (Figures S7C and S7D). For most individuals, the cosine score indicated a better fit of the drift-diffusion model to the observed BOLD trajectories (Figure 5F).

In summary, the WM signals in the visual cortex corroborated the behavioral findings, confirming the implications of the drift-diffusion dynamics.

RNNs with drift dynamics reproduce the human data

We demonstrated that both behavioral and cortical responses support the importance of drift dynamics in explaining the temporal evolution of biases, deriving their implications from phenomenological models. We utilized task-optimized RNNs^{4,49,50} to investigate whether simple network mechanisms

late-DM condition (Figure 5A, right), resulting in a transient drop in stimulus-specific bias (Figure 5B, right) around the discrimination epoch. To properly compare BOLD signals with model predictions, we addressed these issues through event-related analysis⁴⁸: we convolved the predicted trajectories of WM with the canonical hemodynamic response function (HRF) and incorporating memory attraction to reference orientation into the model prediction through regression weights (see STAR Methods). With these corrections, we examined whether cortical signals align with the implications of the drift-diffusion dynamics on decision-consistent bias.

The decision-consistent bias in BOLD signals was estimated as follows: (1) mean decoding errors conditioned on choice at each time point were calculated for near-reference trials (Figure 5C, colored shades), (2) piecewise linear functions were fitted to these error trajectories (Figure 5C, dashed lines), and (3) the bias was quantified by averaging the deviations of the fitted functions from zero at the estimation epoch. Consistent with diffusion-only and drift-diffusion dynamics, we found that the bias increased significantly with decision timing (Figure 5D). Importantly, the decision-timing-dependent measure of bias increase derived from the inferred trajectories

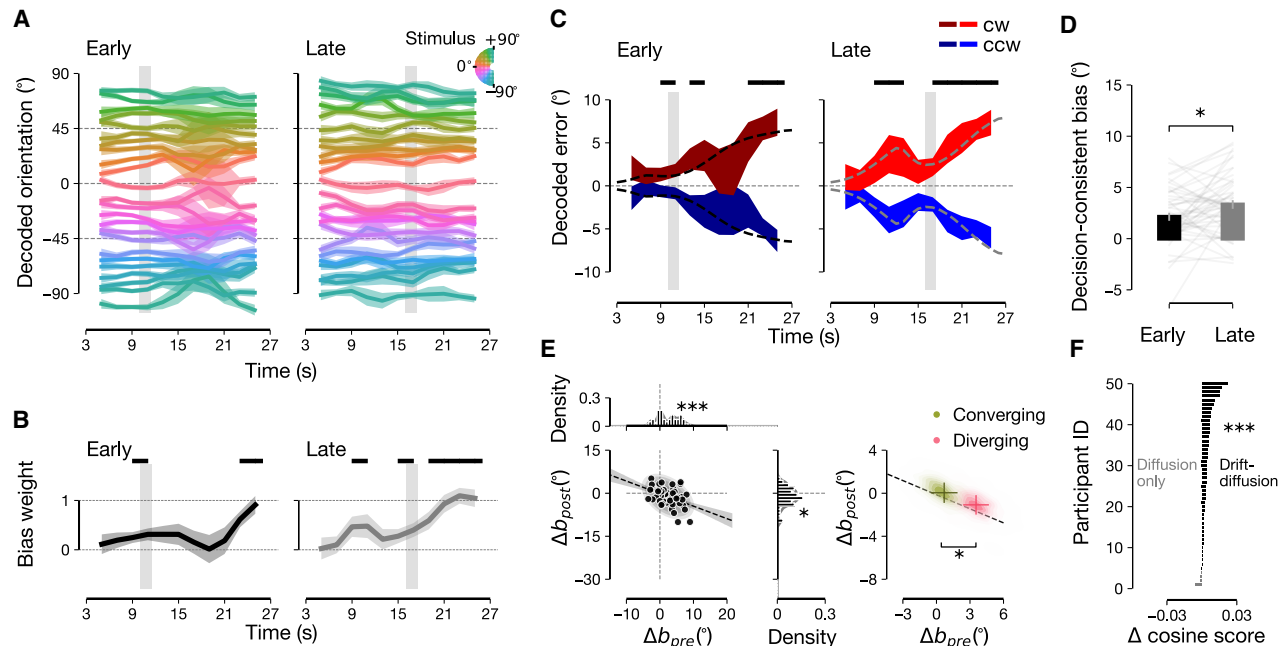


Figure 5. Cortical signals of stimulus-specific drift

(A and B) Evolution of stimulus-specific bias, depicted by decoded orientations (A) and bias weight relative to behavioral estimation bias (B). (C) Evolution of decision-consistent bias, depicted by choice-conditioned decoded errors. Dashed lines mark the BOLD dynamics constructed from the piecewise linear fit (see STAR Methods). (A–C) Gray bars indicate decision timing, with a 4-s hemodynamic delay. Shades, \pm SEMs across trials. (B and C) Black bars mark significant non-zero bias points (B) or between-condition differences (C) $p < 0.05$, permutation test, Bonferroni-corrected for time points. (D) Decision-consistent biases at early and late DM, estimated from the model fit described in (C): gray lines, individuals ($p = 0.0176$, paired t test). (E) Decision-timing-dependent changes in the pre-decision and post-decision biases, estimated from the model fit described in (C) (format as in Figure 4F; one-sample t test, Δb_{pre} , $p < 10^{-5}$; Δb_{post} , $p = 0.0142$; correlation between Δb_{pre} and Δb_{post} , $r = -0.431$, $p = 0.0018$; distance along regression line, 3.069° , $p = 0.0451$). (F) Model comparison in cosine scores of decoded cortical signals (paired t test, $p < 10^{-6}$). *** $p < 0.001$, * $p < 0.05$.

could underlie WM and DM interaction beyond the phenomenological level.

We trained 50 independent RNNs using a task equivalent to that for humans and a joint loss function that penalizes both discrimination and estimation errors (Figure 6A). Given the spatial separation of the stimulus and reference, we fed these inputs to distinct RNN populations. The stimulus inputs were assumed to have greater variability near oblique compared with cardinal orientations, consistent with the efficient coding principle.^{27,36} This heterogeneity in variability prompted the RNNs to drift toward oblique orientations, as this reduced the overall training loss.

The trained RNNs exhibited all the features characteristic of human data, as implied by the drift-diffusion model. Stimulus-specific bias increased over time (Figures 6B, 6E, and 6F), and decision-consistent bias grew (Figures 6C and 6G) with decision timing. A negative correlation was found between Δb_{post} and Δb_{pre} (Figure 6H, left), while Δb_{pre} and Δb_{post} were amplified in orientations with diverging drift (Figure 6H, right).

The RNNs displayed sensory drive effects linked to the reference presentation during the discrimination epoch (transient dips in Figures 6E–6G), as assumed in our BOLD-response model. This is because all neurons in the RNNs, including those receiving the reference, contribute to the estimation (Figure 6A). The RNNs

also exhibited increased variability in estimation errors in near-reference trials, indicating choice-induced bias, which increased with decision timing (Figure 6D, solid symbols). Importantly, these increases in estimation errors in near-reference trials vanished when RNNs were penalized only for estimation errors (Figure 6D, dashed symbols), highlighting DM's critical role in generating choice-induced bias.

In summary, training RNNs to minimize both discrimination and estimation errors while imposing drift dynamics is sufficient for RNNs to display the main features of the estimation biases seen in human data.

RNN mechanism for decision-formation and choice-induced bias

The task-optimized RNNs exhibit a characteristic of choice-induced bias (Figure 6D), mirroring human behavior (Figures 1I, 3C, S1A, and S1B). Investigating the formation of DM from WM and its impact on WM in RNNs may reveal the neural mechanisms underlying this bias, which remains elusive.

To avoid complications from drift dynamics, we studied the interplay between DM and WM in “homogeneous RNNs” trained on inputs devoid of orientation-specific variability, which therefore do not display drift. The average connectivity matrix showed a block-wise structure reflecting separate input and output

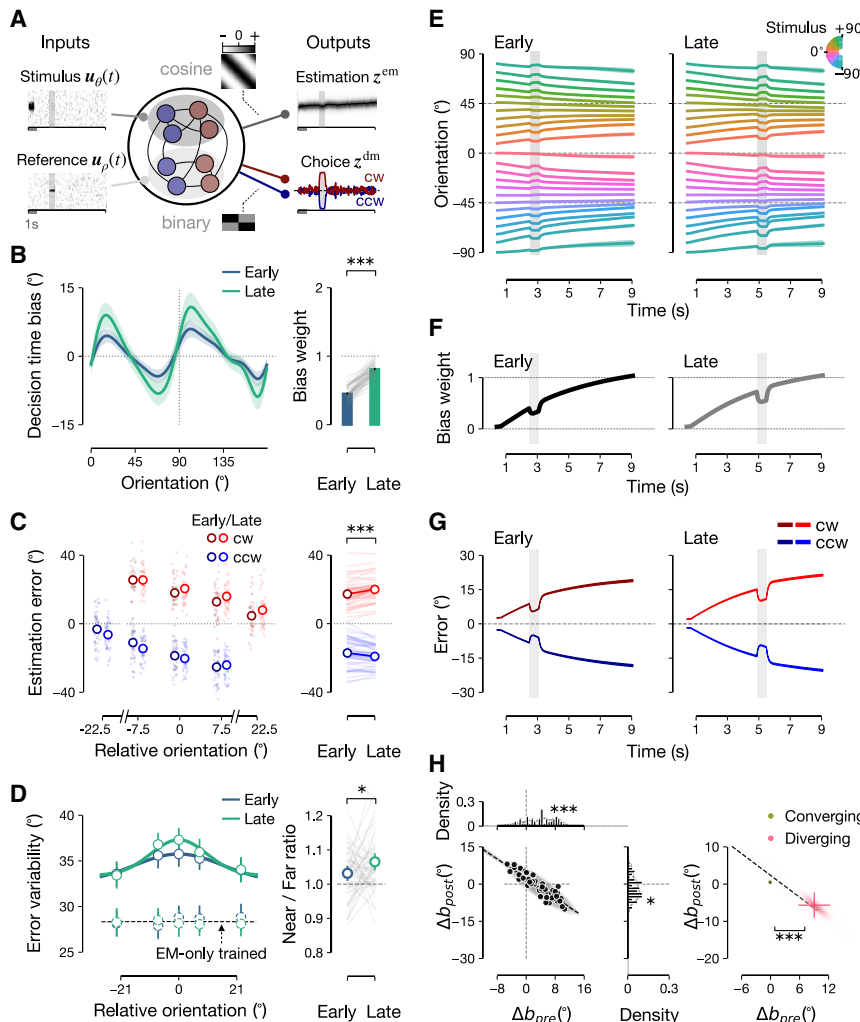


Figure 6. Biases in task-optimized RNNs

(A) RNN architecture with time-varying inputs and outputs. (B–D) Stimulus-specific bias growth (B), decision-timing-dependent decision-consistent biases (C), and near-reference variability (D). Format as in Figures 3A–3C. In (D), solid and dashed lines represent the original and estimation-loss-only RNNs. Paired t test, $p < 10^{-10}$ in (B); $p < 10^{-10}$ in (C); Wilcoxon signed-rank test, $p = 0.0390$ in (D). (E–G) Evolution of stimulus-specific biases (E and F) and decision-consistent biases (G). Format as in Figures 5A–5C, without hemodynamic convolution. (H) Decision-timing-dependent changes in pre-decision and post-decision biases. Format as in Figure 4F. One-sample t test, Δb_{pre} , $p < 10^{-7}$; Δb_{post} , $p = 0.0277$; correlation between Δb_{pre} and Δb_{post} , $r = -0.897$, $p < 10^{-10}$; distance along regression line, 11.491° , $p < 10^{-4}$; ** $p < 0.01$, *** $p < 0.001$.

The DM mechanism can be analyzed geometrically through state-space analysis (Figures 7D–7I). In a 2D-principal-component analysis (PCA) space defined by \mathbf{r}_θ (see STAR Methods), the memory manifolds of the three populations form a ring, initially with \mathbf{r}_ρ^{CW} and \mathbf{r}_ρ^{CCW} being rotated about 45° in the opposite directions from \mathbf{r}_θ (Figures 7D–7F, dotted circles). During DM, reference input vectors (\mathbf{I}_ρ^{ext} , Figure 7G) are added to \mathbf{r}_ρ , making \mathbf{r}_ρ^{CW} (winning population) expand and \mathbf{r}_ρ^{CCW} (losing population) contract (Figures 7D and 7E, arrows on circles). In geometrical terms, correct DM is achieved through a “rotation-addition” mechanism.

pathways (Figure 7A), leading us to analyze interactions among three subpopulations: units receiving stimulus input (\mathbf{r}_θ), those receiving reference input and favoring the CW (\mathbf{r}_ρ^{CW}) and CCW choices (\mathbf{r}_ρ^{CCW}) (Figure 7B).

All three populations exhibit bump-like activity through “feed-forward” (\mathbf{r}_θ to \mathbf{r}_ρ) and “feedback” (\mathbf{r}_ρ to \mathbf{r}_θ) connections. The peaks of \mathbf{r}_θ encode and maintain stimulus orientations faithfully, while \mathbf{r}_ρ^{CW} and \mathbf{r}_ρ^{CCW} shift in CCW and CW directions, respectively (Figure 7C). These shifts reflect \mathbf{r}_θ -to- \mathbf{r}_ρ connectivity, well-approximated by scaled rotations in opposite directions (Figure 7B). Asymmetric connections and shifted memory representations resemble a head-direction system where interactions between opposing populations update direction in response to velocity signals.⁵¹ Similarly, the input from \mathbf{r}_ρ updates \mathbf{r}_θ with a reference input (Figure 7C, shaded horizontal bars).

The feedforward dynamics from \mathbf{r}_θ to \mathbf{r}_ρ underlies DM. Suppose the stimulus is at 0° , and reference is at -7.5° , placing the stimulus CW to the reference. Reference onset increases \mathbf{r}_ρ^{CW} , already rotated CCW, while \mathbf{r}_ρ^{CCW} decreases (Figure 7C, bottom). The DM-mapping matrix reads this amplitude asymmetry into a choice that the stimulus is CW to the reference.

During DM, feedback dynamics implement choice-induced bias. Before DM, the feedback from \mathbf{r}_ρ^{CW} and \mathbf{r}_ρ^{CCW} to \mathbf{r}_θ is balanced, keeping memory at the stimulus orientation. However, with \mathbf{I}_ρ^{ext} on, feedback from the winning population dominates, updating \mathbf{r}_θ in the choice-consistent direction, inducing bias (Figure 7F). After \mathbf{I}_ρ^{ext} is off, \mathbf{r}_ρ quickly returns toward its pre-DM state (Figures S8I and S8J) and ceases reference-related influence. This supports our phenomenological model’s assumption (Figures 1C, 1D, 2A, and 2E): choice induces an immediate, pulse-like update of memory states only during the discrimination epoch.

Further analysis revealed that choice-induced bias ultimately arises from a displacement in feedback dynamics. By linearizing the dynamics along the memory manifold, we found that the feedback rotation opposes the feedforward rotation but over-rotates with a displacement (denoted by φ ; gray area in Figure 7H for over-rotation in $\mathbf{r}_\rho^{CW} \rightarrow \mathbf{r}_\theta$). This creates an imbalance in feedback inputs $\mathbf{I}_{\rho \rightarrow \theta}^{CW}$ and $\mathbf{I}_{\rho \rightarrow \theta}^{CCW}$ with \mathbf{I}_ρ^{ext} on, causing a choice-induced bias proportional to both \mathbf{I}_ρ^{ext} and $\sin(\varphi)$ (Figures 7I and S8C–S8E; Method S4.2). The essential role of the feedback connections is evident when they are ablated before training: discrimination remains intact, but the choice-induced bias does not emerge

(Figures S8K–S8M). Furthermore, decision variable strength scales with choice-induced bias magnitude (Figure S8F), underscoring its functional significance in enabling robust DM under noise.

RNN mechanism of the interaction between stimulus-specific and choice-induced biases

Building on our characterization of how the homogeneous RNNs instantiate choice-induced bias, we revisited the original heterogeneous RNNs—those exhibiting stimulus-specific bias—to identify the mechanism mediating the interaction between stimulus-specific and choice-induced biases. Despite differences in recurrent connectivity (Figure S8B), both RNNs represent stimulus orientations with ring manifolds in similar low-dimensional subspaces. This allows us to project the heterogeneous RNN responses onto the homogeneous RNN state space (Figures 8A–8C).

In the heterogeneous RNNs, ring manifolds in r_ρ were warped into elliptic shapes before the reference input (Figures 8A and 8B), with cardinal orientations more sparsely represented than oblique ones, consistent with efficient coding theories for sensory network.³² In drift dynamics, cardinal and oblique orientations correspond to diverging and converging stimuli, respectively. The warped geometry in r_ρ effectively reverts to a circle shape in r_θ by elongating along the minor axis more than the major axis (Figure 8C; Method S4.1). This anisotropic elongation causes I_ρ^{ext} to have a stronger influence at diverging orientations (Figure 8D) than converging orientations (Figure 8E), resulting in greater choice-induced bias for diverging orientations (Figures 8C and 8D, dotted circles). Stimulus-specific drift further amplifies this effect: perturbations near diverging and converging stimuli are magnified and mitigated, respectively (Figures 4B and 4C), leading to larger biases for diverging orientations (Figure 8F). Human behavioral data validated this prediction: errors in near-reference trials were most biased around cardinal orientations and declined toward oblique orientations (Figure 8G).

In conclusion, the warping geometry of orientation representation and its anisotropic elongation are the key mechanisms mediating the intricate interplay of stimulus-specific, choice-induced, and decision-consistent biases in RNNs.

DISCUSSION

Two unique aspects of our paradigm enable us to identify decision-steered attractor dynamics as a source from which two crucial biases, stimulus-specific and decision-consistent biases, unfold interactively. First, the prolonged delay allows us to probe memory states at sufficiently distant moments through behavioral and neural measurements. Second, the mnemonic discrimination task eliminates sensory access to the target stimulus and thus prevents decision processes from interfering with sensory encoding, unlike previous studies.^{2,8–10,12} With this paradigm, we demonstrated that stimulus-specific bias intensifies over the delay while guiding decisions, and decision-consistent bias increases with decision timing in a stimulus-specific manner.

Analyzing the dynamics of task-optimized RNNs offered valuable insights into the WM-DM interactions. Simplification into three subpopulations, one for WM and two for DM, demon-

strated how categorical decisions emerge from continuous orientation memory and how decisions immediately update memory, creating choice-induced bias. Central to both processes was asymmetric connectivity among these subpopulations, modeled as opposing rotation matrices, with their degree and scaling determining bias strength. This network property, within attractor dynamics, predicted stronger choice-induced bias for orientations with diverging drift, which was confirmed by behavioral data. Further, targeted ablation revealed that feedback from DM to EM populations, which causes the bias, can enhance decision robustness against noise, offering a rational basis for its presence.

Drift dynamics are not the sole source of stimulus-specific bias. Both our phenomenological and RNN models assume that sensory encoding variability—grounded in efficient encoding^{27,28}—also contributes.^{36,37} While downstream readout computations play a role as well,^{27–29,52} the observed growth of stimulus-specific bias in behavioral, BOLD, and RNN data highlights ongoing WM updating via stimulus-specific drift during the delay. Notably, our use of “stimulus-specific drift” differs from “memory drift” in prior literature, which typically refers to random shifts in bump activity within a trial.^{14,53} Many prior studies attribute such shifts to noise-driven diffusion.^{13,30} By contrast, stimulus-specific drift refers to a systematic drift toward fixed attractors amid diffusion, evident when averaged across trials. For that matter, our BOLD-based analyses provide the first neural evidence of stimulus-specific bias growth through systematic drift over tens of seconds.

Stimulus-specific drift toward oblique orientations prompts questions about its mechanisms. Panichello et al.²⁶ demonstrated that discrete attractor dynamics in delayed color estimation could reduce errors by biasing memory toward frequently encountered stimuli. However, as in our and prior studies,^{24,37} orientation estimates are repelled from the frequent, cardinal orientations^{27,28,36} (similarly for location memory³³). Thus, placing attractors around frequent stimuli does not apply to orientation. Instead, our phenomenological and RNN models propose that sensory input to the WM system varies according to the efficient encoding principle, explaining biases and variances inconsistent with traditional attractor models. Recent work³² indicates that orientation error evolution cannot be fully explained by single-module attractor models, emphasizing the roles of sensory and memory network interactions. Further work is needed to understand how these interactions influence error patterns across different features like color and orientation.

Our work clarifies source attributions for decision-consistent bias by differentiating drift toward attractors, stochastic noise,^{13,30} and choice-induced bias.^{2,9} Notably, stimulus-specific drift introduces a dynamic component: it initially biases memory, which biases choices, and then continues to bias memory in line with the biased choice. This has two key implications. First, models of DM and WM should account for how choices feed back into memory. Second, accurately explaining decision-consistent bias requires accounting for drift dynamics alongside stochastic noise, as neglecting this may misattribute post-decision drift to choice-induced bias. Additionally, it offers a perspective on confirmation bias,^{54–56} suggesting it may arise

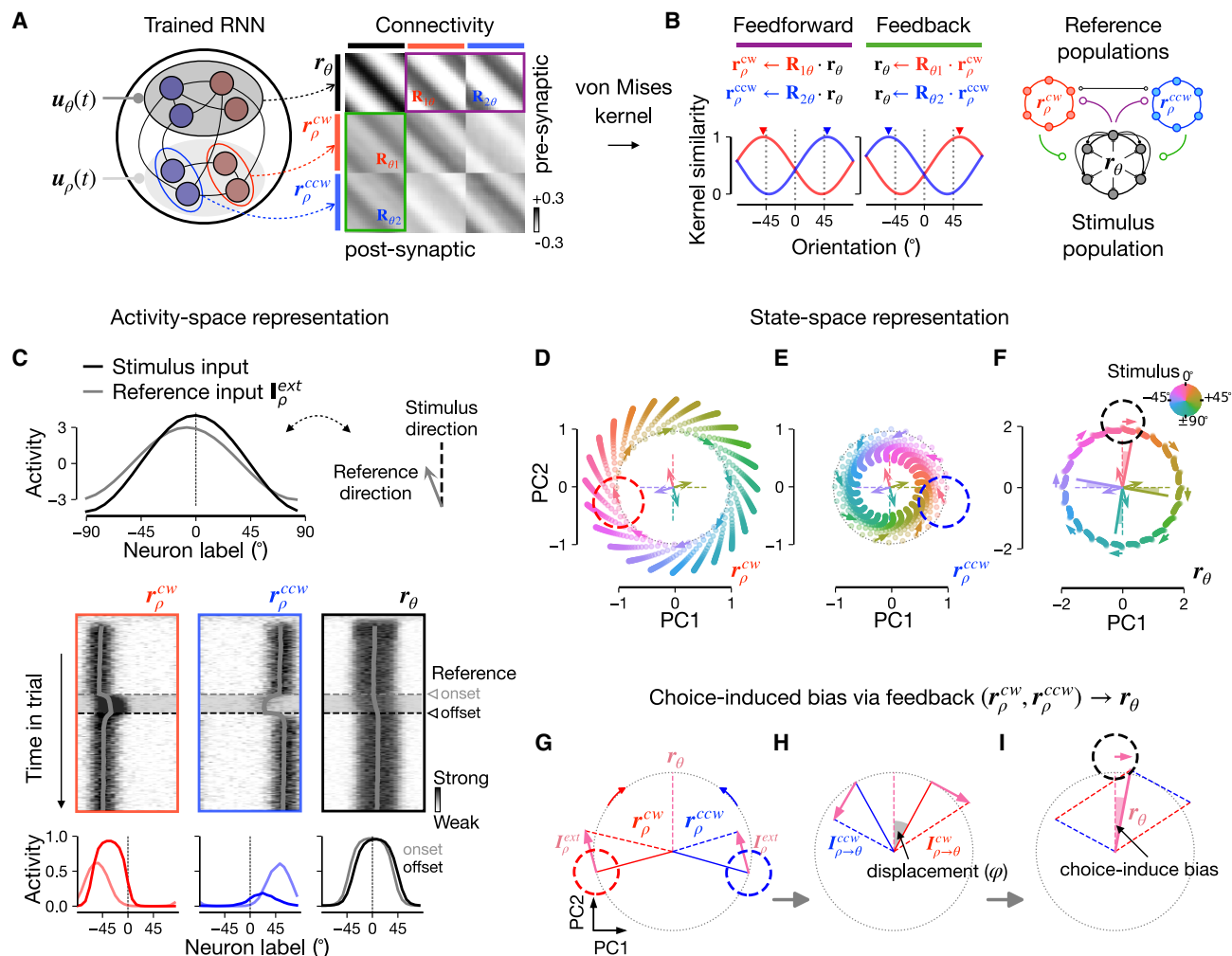


Figure 7. Mechanisms of decision-making and choice-induced bias in homogeneous RNNs

(A) Subpopulations and average connectivities.

(B) Left: scaled-rotation approximation of feedforward and feedback connections: r_θ , stimulus-receiving units; r_ρ^{CW}/r_ρ^{CCW} , reference-receiving, CW/CCW-projecting units. Right: three-ring system with rotation-based recurrent interactions.

(C) Activity changes during decision: top: input profiles over labeled neurons; middle: time courses of r_ρ^{CW} , r_ρ^{CCW} , and r_θ , with discrimination epoch marked with triangles; and bottom: activity snapshots at the onset (light) and offset (dark) of reference.

(D–F) Geometrical analysis of winning (r_ρ^{CW} , D) and losing (r_ρ^{CCW} , E) reference units and stimulus units (r_θ , F) in 2D state space of r_θ during discrimination: isotropic rings, initial memory states; color saturation, time tracked for different stimulus orientations; short arrows, rotation directions; dashed lines and radial arrows, stimulus and reference input for four sample orientations. Dashed circles spotlight the rotation dynamics for 0° stimulus.

(G–I) Linear description of choice-induced bias using low-rank approximation of trained J (see Method S4.1). As r_ρ^{CW} and r_ρ^{CCW} gravitate toward reference input (r_ρ^{ext}), they shift outward and inward, respectively (G). Their feedback to r_θ rotates by a displacement φ (H), yielding summed inputs of r_ρ^{CW} and r_ρ^{CCW} that bias activity in choice-consistent direction (black dashed circle in I).

from decision-consistent bias carried from pre-decision to post-decision phases, facilitated by stimulus-specific drift.

Beyond drift-diffusion dynamics, we provided a mechanistic account of choice-induced bias. Our phenomenological models assumed that decision-formation transiently shifts memory states in the chosen direction during DM, supported by the RNNs. While this bias is attributed to the DM epoch, it may also originate earlier in encoding or later in decoding, with their contributions varying by task structure. Previous accounts^{2,10} suggested that choice-induced bias arises from a non-uniform weighting strategy optimized for DM and subsequently reused

for estimation. This reliance on readout optimization makes it difficult to learn a stable decision boundary when reference inputs vary across trials. Consistent references may allow for late-stage readout strategies via selective information flow⁵⁷ or memory recall.¹² In our paradigm, these optimization or selection strategies seem unlikely, as the reference varies each trial and is briefly available during DM. Further, unlike earlier proposals,^{9,12} situating choice-induced bias during DM predicts specific neural trajectories, verifiable by examining the decision-related and memory-related neural responses with high temporal resolution, as shown by our RNNs.

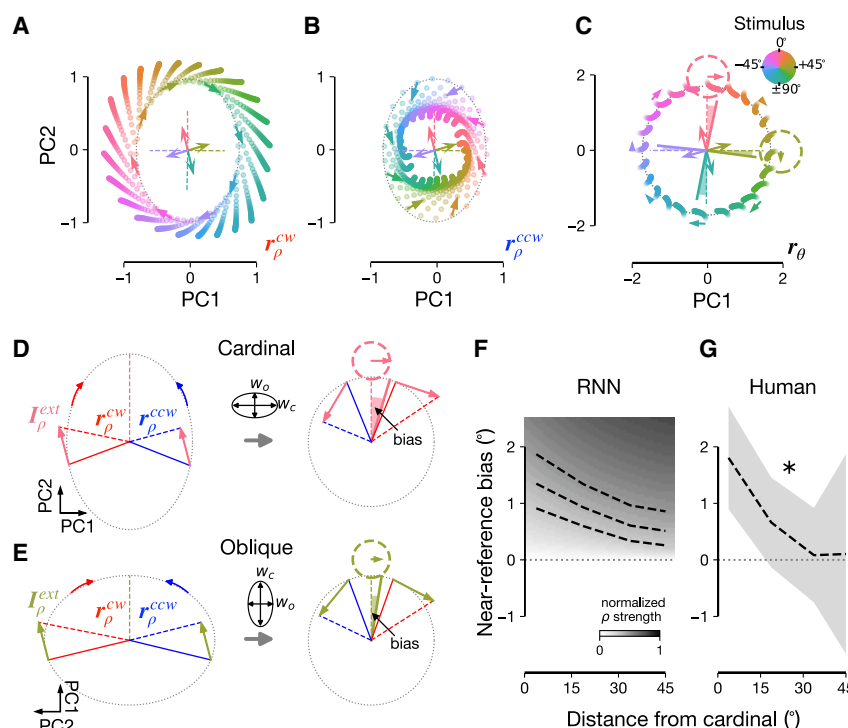


Figure 8. Orientation-dependent biases in heterogeneous RNNs

(A–C) Drift-warped geometry of r_{ρ}^{CW} , r_{ρ}^{CCW} , and r_{θ} during discrimination, shown in the same 2D state space and format as Figures 7D–7F.

(D and E) Linear description of choice-induced bias for cardinal (D, 0°) and oblique (E, 45°) stimuli. The feedback from r_{ρ} to r_{θ} is anisotropic, producing greater elongation along cardinal (w_c) than oblique (w_o) orientations (ellipse above the gray arrows). With $w_c > w_o$, choice-induced bias is more pronounced near cardinal orientations (right panels in D and E; see Method S4.2).

(F) Orientation dependency of near-reference bias. As the reference signal increases in strength, near-reference bias increases while maintaining consistent orientation dependence (darker gray scales with three example bias patterns for different ρ values). (G) Similar orientation-dependent bias observed in human behavior. One-sample t test on individual slopes, $p = 0.0304$. $*p < 0.05$.

Our work offers novel insights into how the brain processes task-relevant features before, during, and after categorical decisions, while optimizing performance in mnemonic discrimination and estimation. These insights warrant further validation and refinement. The mechanism by which our RNNs instantiate choice-induced bias—asymmetric connections and population dynamics—can be explored through synaptic connectivity and state-space dynamics.⁵⁸ Our integrated account of stimulus-specific and decision-consistent biases can also be extended to incorporate effects of memory load^{24,59,60} and serial dependence,^{61–63} which may modulate WM dynamics, possibly through divisive normalization³⁷ or short-term synaptic plasticity.^{64,65} Overall, our work highlights the necessity of considering WM dynamics to fully understand perceptual biases with multiple origins, previously investigated in isolation.

RESOURCE AVAILABILITY

Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by the lead contact, Sang-Hun Lee, at (visionsl@snu.ac.kr).

Materials availability

This study did not produce new materials.

Data and code availability

Raw behavior and fMRI data are publicly available at OpenNeuro: [ds005381](https://doi.org/10.1016/j.neuron.2025.07.003). Processed data are available at Open Science Framework: <https://osf.io/6q95m>. Original Python code for all analyses and figure generation is available at https://github.com/hyunwoogu/dynamic_bias.

ACKNOWLEDGMENTS

We thank Justin L. Gardner, Albert Compte, and Seth W. Egger for helpful comments. We thank Rosanne L. Rademaker, Matthias Fritsche, and Denis Schluppeck for making their data public. S.L. was supported by STI2030-Major Projects, nos. 2021ZD0203700 and 2021ZD0203705. S.L. acknowledges the support of the Shanghai Frontiers Science Center of Artificial Intelligence and Deep Learning and the NYU-ECNU Institute of Brain and Cognitive Science at NYU Shanghai. S.-H.L. was supported by the National Research Foundation of Korea (NRF) funded by the Ministry of Science and Information and Communications Technology (grant nos. NRF-2021R1F1A1052020, NRF-2018R1A4A1025891, and RS-2024-00349515) and by the Korea Basic Science Institute (National Research Facilities and Equipment Center) grant funded by the Ministry of Education (grant no. RS-2024-00435727).

AUTHOR CONTRIBUTIONS

J. Lee, S.K., J. Lim, H.-J.L., M.J.C., D.-g.Y., and S.-H.L. designed research. H. G., S.L., and S.-H.L. performed research. H.G., H.-J.L., H.L., S.K., J. Lim, and J.H.R. contributed unpublished reagents/analytic tools. H.G. wrote the first draft of the paper. H.G., S.L., and S.-H.L. edited and wrote the paper.

DECLARATION OF INTERESTS

The authors declare no competing interests.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- EXPERIMENTAL MODEL AND SUBJECT DETAILS
- METHOD DETAILS
 - Experiments

- Analysis of data
- Phenomenological models: diffusion-only and drift-diffusion models
- Recurrent neural network model
- **QUANTIFICATION AND STATISTICAL PROCEDURES**

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.neuron.2025.07.003>.

Received: November 13, 2024

Revised: May 22, 2025

Accepted: July 3, 2025

REFERENCES

1. Ma, W.J., and Jazayeri, M. (2014). Neural coding of uncertainty and probability. *Annu. Rev. Neurosci.* 37, 205–220. <https://doi.org/10.1146/annurev-neuro-071013-014017>.
2. Jazayeri, M., and Movshon, J.A. (2007). A new perceptual illusion reveals mechanisms of sensory decoding. *Nature* 446, 912–915. <https://doi.org/10.1038/nature05739>.
3. Gold, J.I., and Shadlen, M.N. (2007). The Neural Basis of Decision Making. *Annu. Rev. Neurosci.* 30, 535–574. <https://doi.org/10.1146/annurev-neuro.29.051605.113038>.
4. Mante, V., Sussillo, D., Shenoy, K.V., and Newsome, W.T. (2013). Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* 503, 78–84. <https://doi.org/10.1038/nature12742>.
5. Bogacz, R., Brown, E., Moehlis, J., Holmes, P., and Cohen, J.D. (2006). The physics of optimal decision making: A formal analysis of models of performance in two-alternative forced-choice tasks. *Psychol. Rev.* 113, 700–765. <https://doi.org/10.1037/0033-295X.113.4.700>.
6. Romo, R., Brody, C.D., Hernández, A., and Lemus, L. (1999). Neuronal correlates of parametric working memory in the prefrontal cortex. *Nature* 399, 470–473. <https://doi.org/10.1038/20939>.
7. Machens, C.K., Romo, R., and Brody, C.D. (2005). Flexible Control of Mutual Inhibition: A Neural Model of Two-Interval Discrimination. *Science* 307, 1121–1124. <https://doi.org/10.1126/science.1104171>.
8. Fritsche, M., and de Lange, F.P. (2019). Reference repulsion is not a perceptual illusion. *Cognition* 184, 107–118. <https://doi.org/10.1016/j.cognition.2018.12.010>.
9. Luu, L., and Stocker, A.A. (2018). Post-decision biases reveal a self-consistency principle in perceptual inference. *eLife* 7, e33334. <https://doi.org/10.7554/eLife.33334>.
10. Zamboni, E., Ledgeway, T., McGraw, P.V., and Schluppeck, D. (2016). Do perceptual biases emerge early or late in visual processing? Decision-biases in motion perception. *Proc. Biol. Sci.* 283, 20160263. <https://doi.org/10.1098/rspb.2016.0263>.
11. Stocker, A.A., and Simoncelli, E.P. (2007). A Bayesian Model of Conditioned Perception. *Adv. Neural Inf. Process. Syst.* 2007, 1409–1416.
12. Luu, L., and Stocker, A.A. (2021). Categorical judgments do not modify sensory representations in working memory. *PLOS Comput. Biol.* 17, e1008968. <https://doi.org/10.1371/journal.pcbi.1008968>.
13. Compte, A., Brunel, N., Goldman-Rakic, P.S., and Wang, X.J. (2000). Synaptic Mechanisms and Network Dynamics Underlying Spatial Working Memory in a Cortical Network Model. *Cereb. Cortex* 10, 910–923. <https://doi.org/10.1093/cercor/10.9.910>.
14. Schneegans, S., and Bays, P.M. (2018). Drift in neural population activity causes working memory to deteriorate over time. *J. Neurosci.* 38, 4859–4869. <https://doi.org/10.1523/JNEUROSCI.3440-17.2018>.
15. Murray, J.D., Jaramillo, J., and Wang, X.J. (2017). Working memory and decision-making in a frontoparietal circuit model. *J. Neurosci.* 37, 12167–12186. <https://doi.org/10.1523/JNEUROSCI.0343-17.2017>.
16. Meister, M.L.R., Hennig, J.A., and Huk, A.C. (2013). Signal Multiplexing and Single-Neuron Computations in Lateral Intraparietal Area During Decision-Making. *J. Neurosci.* 33, 2254–2267. <https://doi.org/10.1523/JNEUROSCI.2984-12.2013>.
17. Wang, X.J. (2008). Decision Making in Recurrent Neuronal Circuits. *Neuron* 60, 215–234. <https://doi.org/10.1016/j.neuron.2008.09.034>.
18. Lemus, L., Hernández, A., Luna, R., Zainos, A., Nácher, V., and Romo, R. (2007). Neural correlates of a postponed decision report. *Proc. Natl. Acad. Sci. USA* 104, 17174–17179. <https://doi.org/10.1073/pnas.0707961104>.
19. Huttenlocher, J., Hedges, L.V., and Duncan, S. (1991). Categories and particulars: Prototype effects in estimating spatial location. *Psychol. Rev.* 98, 352–376. <https://doi.org/10.1037/0033-295X.98.3.352>.
20. Bae, G.Y. (2021). Neural evidence for categorical biases in location and orientation representations in a working memory task. *Neuroimage* 240, 118366. <https://doi.org/10.1016/j.neuroimage.2021.118366>.
21. Blake, R., Cepeda, N.J., and Hiris, E. (1997). Memory for Visual Motion. *J. Exp. Psychol. Hum. Percept. Perform.* 23, 353–369. <https://doi.org/10.1037/0096-1523.23.2.353>.
22. Rauber, H.J., and Treue, S. (1998). Reference repulsion when judging the direction of visual motion. *Perception* 27, 393–402. <https://doi.org/10.1068/p270393>.
23. de Gardelle, V., Kouider, S., and Sackur, J. (2010). An oblique illusion modulated by visibility: Non-monotonic sensory integration in orientation processing. *J. Vis.* 10, 6. <https://doi.org/10.1167/10.10.6>.
24. Pratte, M.S., Park, Y.E., Rademaker, R.L., and Tong, F. (2017). Accounting for stimulus-specific variation in precision reveals a discrete capacity limit in visual working memory. *J. Exp. Psychol. Hum. Percept. Perform.* 43, 6–17. <https://doi.org/10.1037/xhp0000302>.
25. Hardman, K.O., Vergauwe, E., and Ricker, T.J. (2017). Categorical working memory representations are used in delayed estimation of continuous colors. *J. Exp. Psychol. Hum. Percept. Perform.* 43, 30–54. <https://doi.org/10.1037/xhp0000290>.
26. Panichello, M.F., DePasquale, B., Pillow, J.W., and Buschman, T.J. (2019). Error-correcting dynamics in visual working memory. *Nat. Commun.* 10, 3366. <https://doi.org/10.1038/s41467-019-11298-3>.
27. Wei, X.X., and Stocker, A.A. (2015). A Bayesian observer model constrained by efficient coding can explain “anti-Bayesian” percepts. *Nat. Neurosci.* 18, 1509–1517. <https://doi.org/10.1038/nn.4105>.
28. Wei, X.X., and Stocker, A.A. (2017). Lawful relation between perceptual bias and discriminability. *Proc. Natl. Acad. Sci. USA* 114, 10244–10249. <https://doi.org/10.1073/pnas.1619153114>.
29. Hahn, M., and Wei, X.X. (2024). A unifying theory explains seemingly contradictory biases in perceptual estimation. *Nat. Neurosci.* 27, 793–804. <https://doi.org/10.1038/s41593-024-01574-x>.
30. Wimmer, K., Nykamp, D.Q., Constantinidis, C., and Compte, A. (2014). Bump attractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nat. Neurosci.* 17, 431–439. <https://doi.org/10.1038/nn.3645>.
31. Eissa, T.L., and Kilpatrick, Z.P. (2023). Learning efficient representations of environmental priors in working memory. *PLOS Comput. Biol.* 19, e1011622. <https://doi.org/10.1371/journal.pcbi.1011622>.
32. Yang, J., Zhang, H., and Lim, S. (2024). Sensory-memory interactions via modular structure explain errors in visual working memory. *eLife* 13, RP95160. <https://doi.org/10.7554/eLife.95160>.
33. Stein, H., Barbosa, J., Rosa-Justicia, M., Prades, L., Morató, A., Galan-Gadea, A., Ariño, H., Martínez-Hernández, E., Castro-Fornieles, J., Dalmau, J., et al. (2020). Reduced serial dependence suggests deficits in synaptic potentiation in anti-NMDAR encephalitis and schizophrenia. *Nat. Commun.* 11, 4250. <https://doi.org/10.1038/s41467-020-18033-3>.
34. Tomassini, A., Morgan, M.J., and Solomon, J.A. (2010). Orientation uncertainty reduces perceived obliquity. *Vision Res.* 50, 541–547. <https://doi.org/10.1016/j.visres.2009.12.005>.

35. Yu, Q., Panichello, M.F., Cai, Y., Postle, B.R., and Buschman, T.J. (2020). Delay-period activity in frontal, parietal, and occipital cortex tracks noise and biases in visual working memory. *PLOS Biol.* 18, e3000854. <https://doi.org/10.1371/journal.pbio.3000854>.
36. Girshick, A.R., Landy, M.S., and Simoncelli, E.P. (2011). Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nat. Neurosci.* 14, 926–932. <https://doi.org/10.1038/nn.2831>.
37. Taylor, R., and Bays, P.M. (2018). Efficient coding in visual working memory accounts for stimulus-specific variations in recall. *J. Neurosci.* 38, 7132–7142. <https://doi.org/10.1523/JNEUROSCI.1018-18.2018>.
38. Ganguli, D., and Simoncelli, E.P. (2014). Efficient Sensory Encoding and Bayesian Inference with Heterogeneous Neural Populations. *Neural Comput.* 26, 2103–2134. https://doi.org/10.1162/NECO_a_00638.
39. Benjamin, A.S., Zhang, L.Q., Qiu, C., Stocker, A.A., and Kording, K.P. (2022). Efficient neural codes naturally emerge through gradient descent learning. *Nat. Commun.* 13, 7972. <https://doi.org/10.1038/s41467-022-35659-7>.
40. Mao, J., and Stocker, A.A. (2024). Sensory perception is a holistic inference process. *Psychol. Rev.* 131, 858–890. <https://doi.org/10.1037/rev0000457>.
41. Morais, M.J., and Pillow, J.W. (2018). Power-law efficient neural codes provide general link between perceptual bias and discriminability. *Adv. Neural Inf. Process. Syst.* 31, 5071–5080.
42. Brouwer, G.J., and Heeger, D.J. (2009). Decoding and Reconstructing Color from Responses in Human Visual Cortex. *J. Neurosci.* 29, 13992–14003. <https://doi.org/10.1523/JNEUROSCI.3577-09.2009>.
43. Rademaker, R.L., Chunharas, C., and Serences, J.T. (2019). Coexisting representations of sensory and mnemonic information in human visual cortex. *Nat. Neurosci.* 22, 1336–1344. <https://doi.org/10.1038/s41593-019-0428-x>.
44. Harrison, W.J., Bays, P.M., and Rideaux, R. (2023). Neural tuning instantiates prior expectations in the human visual system. *Nat. Commun.* 14, 5320. <https://doi.org/10.1038/s41467-023-41027-w>.
45. Serences, J.T., Ester, E.F., Vogel, E.K., and Awh, E. (2009). Stimulus-specific delay activity in human primary visual cortex. *Psychol. Sci.* 20, 207–214. <https://doi.org/10.1111/j.1467-9280.2009.02276.x>.
46. Master, S.L., Li, S., and Curtis, C.E. (2024). Trying Harder: How Cognitive Effort Sculpt Neural Representations during Working Memory. *J. Neurosci.* 44, e0060242024. <https://doi.org/10.1523/JNEUROSCI.0060-24.2024>.
47. Kriegeskorte, N., Mur, M., and Bandettini, P. (2008). Representational similarity analysis – connecting the branches of systems neuroscience. *Front. Syst. Neurosci.* 2, 4. <https://doi.org/10.3389/neuro.06.004.2008>.
48. Friston, K.J., Fletcher, P., Josephs, O., Holmes, A., Rugg, M.D., and Turner, R. (1998). Event-related fMRI: Characterizing differential responses. *Neuroimage* 7, 30–40. <https://doi.org/10.1006/nimg.1997.0306>.
49. Yang, G.R., Joglekar, M.R., Song, H.F., Newsome, W.T., and Wang, X.J. (2019). Task representations in neural networks trained to perform many cognitive tasks. *Nat. Neurosci.* 22, 297–306. <https://doi.org/10.1038/s41593-018-0310-2>.
50. Dubreuil, A., Valente, A., Beiran, M., Mastrogiuseppe, F., and Ostojic, S. (2022). The role of population structure in computations through neural dynamics. *Nat. Neurosci.* 25, 783–794. <https://doi.org/10.1038/s41593-022-01088-4>.
51. Xie, X., Hahnloser, R.H.R., and Seung, H.S. (2002). Selectively grouping neurons in recurrent networks of lateral inhibition. *Neural Comput.* 14, 2627–2646. <https://doi.org/10.1162/089976602760408008>.
52. Birman, D., and Gardner, J.L. (2019). A flexible readout mechanism of human sensory representations. *Nat. Commun.* 10, 3500. <https://doi.org/10.1038/s41467-019-11448-7>.
53. Wolff, M.J., Jochim, J., Akyürek, E.G., Buschman, T.J., and Stokes, M.G. (2020). Drifting codes within a stable coding scheme for working memory. *PLOS Biol.* 18, e3000625. <https://doi.org/10.1371/journal.pbio.3000625>.
54. Talluri, B.C., Urai, A.E., Tsetsos, K., Usher, M., and Donner, T.H. (2018). Confirmation Bias through Selective Overweighting of Choice-Consistent Evidence. *Curr. Biol.* 28, 3128–3135.e8. <https://doi.org/10.1016/j.cub.2018.07.052>.
55. Glickman, M., Moran, R., and Usher, M. (2022). Evidence integration and decision confidence are modulated by stimulus consistency. *Nat. Hum. Behav.* 6, 988–999. <https://doi.org/10.1038/s41562-022-01318-6>.
56. Lange, R.D., Chatteraj, A., Beck, J.M., Yates, J.L., and Haefner, R.M. (2021). A confirmation bias in perceptual decision-making due to hierarchical approximate inference. *PLOS Comput. Biol.* 17, e1009517. <https://doi.org/10.1371/journal.pcbi.1009517>.
57. Yang, G.R., Murray, J.D., and Wang, X.J. (2016). A dendritic disinhibitory circuit mechanism for pathway-specific gating. *Nat. Commun.* 7, 12815. <https://doi.org/10.1038/ncomms12815>.
58. Paninski, L., and Cunningham, J.P. (2018). Neural data science: accelerating the experiment-analysis-theory cycle in large-scale neuroscience. *Curr. Opin. Neurobiol.* 50, 232–241. <https://doi.org/10.1016/j.conb.2018.04.007>.
59. van den Berg, R., Shin, H., Chou, W.-C., George, R., and Ma, W.J. (2012). Variability in encoding precision accounts for visual short-term memory limitations. *Proc. Natl. Acad. Sci. USA* 109, 8780–8785. <https://doi.org/10.1073/pnas.1117465109>.
60. Bays, P.M. (2014). Noise in neural populations accounts for errors in working memory. *J. Neurosci.* 34, 3632–3645. <https://doi.org/10.1523/JNEUROSCI.3204-13.2014>.
61. Fritzsche, M., Mostert, P., and de Lange, F.P. (2017). Opposite Effects of Recent History on Perception and Decision. *Curr. Biol.* 27, 590–595. <https://doi.org/10.1016/j.cub.2017.01.006>.
62. Bliss, D.P., Sun, J.J., and D’Esposito, M. (2017). Serial dependence is absent at the time of perception but increases in visual working memory. *Sci. Rep.* 7, 14739. <https://doi.org/10.1038/s41598-017-15199-7>.
63. Markov, Y.A., Tiurina, N.A., and Pascucci, D. (2024). Serial dependence: A matter of memory load. *Heliyon* 10, e33977. <https://doi.org/10.1016/j.heliyon.2024.e33977>.
64. Kilpatrick, Z.P. (2018). Synaptic mechanisms of interference in working memory. *Sci. Rep.* 8, 7879. <https://doi.org/10.1038/s41598-018-25958-9>.
65. Barbosa, J., Stein, H., Martinez, R.L., Galan-Gadea, A., Li, S., Dalmau, J., Adam, K.C.S., Valls-Solé, J., Constantinidis, C., and Compte, A. (2020). Interplay between persistent activity and activity-silent dynamics in the prefrontal cortex underlies serial biases in working memory. *Nat. Neurosci.* 23, 1016–1024. <https://doi.org/10.1038/s41593-020-0644-4>.
66. Esteban, O., Markiewicz, C.J., Blair, R.W., Moodie, C.A., Isik, A.I., Erramuzpe, A., Kent, J.D., Goncalves, M., DuPre, E., Snyder, M., et al. (2019). fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat. Methods* 16, 111–116. <https://doi.org/10.1038/s41592-018-0235-4>.
67. Gardner, J.L., Merriam, E.P., Schluppeck, D., and Larsson, J. (2018). MGL: Visual psychophysics stimuli and experimental design package, [Version 2.0]. Zenodo. <https://doi.org/10.5281/zenodo.1299497>.
68. Gardner, J.L., Merriam, E.P., Schluppeck, D., Besle, J., and Heeger, D.J. (2018). mrTools: Analysis and visualization package for functional magnetic resonance imaging data, [Version 4.7]. Zenodo. <https://doi.org/10.5281/zenodo.1299483>.
69. Engel, S.A., Rumelhart, D.E., Wandell, B.A., Lee, A.T., Glover, G.H., Chichilnisky, E.J., and Shadlen, M.N. (1994). fMRI of human visual cortex. *Nature* 369, 525. <https://doi.org/10.1038/369525a0>.
70. Choe, K.W., Blake, R., and Lee, S.-H. (2014). Dissociation between Neural Signatures of Stimulus and Choice in Population Activity of Human V1 during Perceptual Decision-Making. *J. Neurosci.* 34, 2725–2743. <https://doi.org/10.1523/JNEUROSCI.1606-13.2014>.
71. Ryu, J., and Lee, S.-H. (2018). Stimulus-Tuned Structure of Correlated fMRI Activity in Human Visual Cortex. *Cereb. Cortex* 28, 693–712. <https://doi.org/10.1093/cercor/bhw411>.

72. Wang, L., Mruczek, R.E.B., Arcaro, M.J., and Kastner, S. (2015). Probabilistic maps of visual topography in human cortex. *Cereb. Cortex* 25, 3911–3931. <https://doi.org/10.1093/cercor/bhu277>.
73. Glasser, M.F., Coalson, T.S., Robinson, E.C., Hacker, C.D., Harwell, J., Yacoub, E., Ugurbil, K., Andersson, J., Beckmann, C.F., Jenkinson, M., et al. (2016). A multi-modal parcellation of human cerebral cortex. *Nature* 536, 171–178. <https://doi.org/10.1038/nature18933>.
74. Lindquist, M.A., Geuter, S., Wager, T.D., and Caffo, B.S. (2019). Modular preprocessing pipelines can reintroduce artifacts into fMRI data. *Hum. Brain Mapp.* 40, 2358–2376. <https://doi.org/10.1002/hbm.24528>.
75. Rademaker, R.L., Bloem, I.M., De Weerd, P., and Sack, A.T. (2015). The impact of interference on short-term memory for visual orientation. *J. Exp. Psychol. Hum. Percept. Perform.* 41, 1650–1665. <https://doi.org/10.1037/xhp0000110>.
76. Green, D.M., and Swets, J.A. (1966). *Signal Detection Theory and Psychophysics* (John Wiley & Sons).
77. Glover, G.H. (1999). Deconvolution of Impulse Response in Event-Related BOLD fMRI. *Neuroimage* 9, 416–429. <https://doi.org/10.1006/nimg.1998.0419>.
78. Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., and Diedrichsen, J. (2016). Reliability of dissimilarity measures for multi-voxel pattern analysis. *Neuroimage* 137, 188–200. <https://doi.org/10.1016/j.neuroimage.2015.12.012>.
79. Pietrini, P., Furey, M.L., Ricciardi, E., Gobbini, M.I., Wu, W.H.C., Cohen, L., Guazzelli, M., and Haxby, J.V. (2004). Beyond sensory images: Object-based representation in the human ventral pathway. *Proc. Natl. Acad. Sci. USA* 101, 5658–5663. <https://doi.org/10.1073/pnas.0400707101>.
80. Risken, H. (1996). *The Fokker-Planck Equation: Methods of Solution and Applications* (Springer). <https://doi.org/10.1007/978-3-642-61544-3>.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
Behavior data (preprocessed)	This paper	Open Science Framework: https://osf.io/6q95m
fMRI data (preprocessed)	This paper	Open Science Framework: https://osf.io/6q95m
Behavior and fMRI data (raw)	This paper	OpenNeuro: https://openneuro.org/datasets/ds005381
Software and algorithms		
fMRIPrep 20.2.0	Esteban et al. ⁶⁶	https://fmriprep.org/en/20.2.0/index.html ; RRID: SCR_016216
MGL	Gardner et al. ⁶⁷	https://github.com/justingardner/mgl , https://zenodo.org/records/1299497
mrTools	Gardner et al. ⁶⁸	https://github.com/justingardner/mrTools , https://zenodo.org/records/1299483
Custom code	This paper	https://github.com/hyunwoogu/dynamic_bias

EXPERIMENTAL MODEL AND SUBJECT DETAILS

This study was conducted in accordance with the guidelines of and under the approval of the Institutional Review Board of Seoul National University. 50 healthy individuals (30 females, 19 – 32 years old; normal or corrected-to-normal vision) completed at least two or all of the three-day sessions of the main fMRI experiment across three days. Each participant provided written informed consent prior to the experiment and was naïve to the purpose of the study.

METHOD DETAILS

Experiments

Experiment stimuli and procedure

Stimuli were generated using MGL⁶⁷ and presented by an LCD projector (60Hz). Participants viewed the stimuli at a visual angle of 22° (width) × 17° (height). The stimuli were displayed within a whole field (radius, 8.5°) gray circular Gaussian envelope aperture on a black background. A black fixation dot (0.07°) and a surrounding black fixation ring (0.83° to 0.9°) were constantly present at the center of the screen. There were 20 runs across the three scanning sessions, each run consisting of 12 trials. Some runs were excluded due to excessive motion (maximum motion across axes in rotation and translation surpassed T2*-weighted voxel size). Participants completed a one-hour practice session a few days before the main fMRI experiment.

Inside the scanner, participants maintained central fixation throughout each run and responded using a button box with linearly aligned keys labeled Key1 to Key4. Before each trial, the fixation dot dilated (0.14°) for 0.5s to cue the stimulus onset. Each trial started with a 1.5s presentation of an alternating (8/3 Hz) donut-shaped oriented grating (spatial frequency, 1 cycle/degree) spanning the peripheral visual field (aperture radii: inner, 2°; outer, 8.5°). The stimulus orientations ranged from 0° to 172.5° with a step size of 7.5°. The target stimulus presentation was followed by a first-epoch delay, a discrimination task, a second-epoch delay, and an estimation task. Two conditions were considered: early DM trials with 4.5s first-epoch and 10.5s second-epoch delays, and late DM trials with 10.5s first-epoch and 4.5s second-epoch delays.

In the discrimination task, participants viewed an oriented reference frame, a virtual line connecting the two yellow nonius dots (mark size, 0.1°) on the fixation ring. The fixation dot turned yellow and transiently dilated to 0.14° for 0.5s to cue the discrimination task onset. Participants indicated whether the target was tilted counter-clockwise or clockwise relative to the reference by pressing the Key2 (CCW) or Key3 (CW) with their left or right thumb, respectively. Relative orientations of the reference to the target stimulus were uniformly selected from [−21°, −4°, 0°, 4°, 21°], whose range approximately matches those of the previous studies.^{2,9} Participants had 1.5s to respond, with their responses recorded without their knowledge, with a 0.5s buffer. If a response was made, the fixation dot dilated to 0.14° and turned blue for 0.75s; if no response was made within 1.5s, the fixation dot dilated to 0.14° and turned red for 0.75s. The reference frame disappeared upon button press.

In the estimation task, the participants reproduced the target stimulus from memory by rotating the two green nonius dots with Key2 (CCW) and Key3 (CW) within the 4.5s. The fixation dot turned green and transiently dilated to eccentricity 0.14° for 0.5s to cue the estimation task onset. Participants confirmed their adjustments by pressing Key1 using their thumbs. If a response was made, the fixation dot dilated to 0.14° and turned blue for 0.75s; if no response was made within 4.5s, the fixation dot dilated to 0.14° and turned red for 0.75s. The starting orientation of the estimation nonius dots was randomly chosen from 0° to 180°. The estimation task was followed by a 5.5s inter-trial interval (ITI). Each trial lasted 28s, with a total run time of 336s. After each run, participants received a summary of their performance.

MRI data acquisition and preprocessing

MR data were collected using a Siemens 3 Tesla Tim Trio with a 32-channel head matrix coil at the Seoul National University Brain Imaging Center. Participants underwent T1-weighted, high-resolution ($0.8 \times 0.8 \times 0.8 \text{ mm}^3$) anatomical scans (repetition time (TR), 2.4s; inversion time (TI), 1s; time to echo (TE), 2.19ms; flip angle (FA), 8°). Over three separate days, they participated in the main T2*-weight fMRI scanning sessions: Day 1 with of retinotopy-mapping run (96s), hemodynamic impulse response function (HIRF) estimation run (96s), and 6 task runs (336s), Day 2 with 8 task runs, and Day 3 with 6 task runs. Scan parameters for the retinotopy-mapping, HIRF, and task runs were: voxel size, $2.3 \times 2.3 \times 2.3 \text{ mm}^3$; TR, 2.0s; TE, 30ms; FA, 77° . After the acquisition of fMRI scanner data, the initial preprocessing steps for the anatomical and functional images followed the fMRIPrep workflow⁶⁶ (version 20.2.0) with field map-free distortion correction option (–use-syn-sdc). For control analyses, we included the projections onto the FSL MNI space (MNI152NLIN6Asym).

ROI definition and voxel selection

The V1, V2, and V3 were defined using standard traveling wave methods.⁶⁹ Two 15° -wide wedge bowties on the vertical and horizontal meridians served as stimuli as in our previous study.⁷⁰ To measure the voxel-wise signal-to-noise ratio (SNR), we used a checkerboard whole-field impulse (radius, 8°) at a 1/24Hz frequency in the HIRF scan. Using the retinotopy scan, subjects' V1, V2, and V3 regions across the left and right hemispheres and across dorsal and ventral areas were defined and combined for BOLD analysis. Voxel-wise SNR was calculated as the stimulus frequency (1/24Hz) amplitude in the HIRF scan divided by the average amplitude of frequencies above the third harmonics, discarding voxels with SNRs under two, following our previous research.⁷¹ For control analyses, we used the publicly available visual atlas for IPS⁷² as well as the frontal regions,⁷³ IFC and DLPFC. Each voxel's time series was converted into percent signal change by dividing by its average over the entire time series. To minimize the artifacts, fMRIPrep-derived confounding variables were regressed out, consisting of white matter, CSF, and six additional three-dimensional motion regressors, along with the discrete cosine transform bases below 0.008Hz to reduce low-frequency components. No additional spatial smoothing was applied. The confounders were regressed out simultaneously to minimize the potential artifacts from the stepwise regression.⁷⁴ Resulting time series were z-scored voxel-by-voxel and run-by-run for further analyses.

Analysis of data

Quantifying the stimulus-specific bias from behavior data

The stimulus-specific bias was quantified from the estimation data or the discrimination data. As for the estimation data, we computed the stimulus-conditioned means of estimation errors $\varepsilon(\theta)$, the difference between the estimation ($\hat{\theta}$) and the stimulus (θ), and fitted a smooth function $\kappa(\theta)$ to $\varepsilon(\theta)$. For each participant, the best-fit smooth function $\hat{\kappa}(\theta) = \mathbf{v}(\theta)^\top \omega^*$ was found by finding ω that minimizes the sum of squared errors,

$$\omega^* = \underset{\omega}{\operatorname{argmin}} \sum_{j=1}^{N_{\text{trial}}} (\mathbf{v}(\theta_j)^\top \omega - \varepsilon(\theta_j))^2 \quad (\text{Equation 1})$$

where $\mathbf{v}(\theta) = [1, v'_1(\theta), \dots, v'_{N_{\text{basis}}}(\theta)]^\top$ with v'_k being the derivative of the von Mises density function with a center of $2j\pi/N_{\text{basis}}$, a pre-cision of $N_{\text{basis}}/2$, and N_{basis} set to 12.

As for the discrimination data, we fitted psychometric functions Ψ to the discrimination choices ($\hat{c}_1, \dots, \hat{c}_{N_{\text{trial}}}$) by maximizing the likelihood of choices $\hat{c} \in \{-1, +1\}$, corresponding to CCW (-1) and CW ($+1$), as follows:

$$L(\hat{c}_1, \dots, \hat{c}_{N_{\text{trial}}}) = \prod_{j=1}^{N_{\text{trial}}} \Psi(\theta_j, \rho_j, \vartheta_j)^{(1+\hat{c}_j)/2} (1 - \Psi(\theta_j, \rho_j, \vartheta_j))^{(1-\hat{c}_j)/2}; \quad (\text{Equation 2})$$

$$\Psi(\theta, \rho, \vartheta) = \lambda + (1 - 2\lambda) \cdot \Phi(\tilde{\rho}; \mu, \sigma^2) \quad (\text{Equation 3})$$

where θ, ρ, ϑ are stimulus orientation, reference orientation, and decision timing ($\vartheta \in \{\text{early}, \text{late}\}$); $\tilde{\rho}$ is the orientation of the reference relative to the stimulus ($\tilde{\rho} = \rho - \theta$); Φ is the cumulative Gaussian distribution function with a mean μ and a standard deviation σ ; λ is the lapse rate. While fitting Ψ , to parameterize the modulation of the stimulus-specific bias by decision timing (Figure 3A), we constrained the stimulus-specific mean of Φ with the best-fit smooth function $\hat{\kappa}(\theta)$, as follows: $\mu = w_\vartheta \hat{\kappa}(\theta)$, where, w_ϑ denotes the bias weight for the early (w_{early}) or late (w_{late}) DM conditions. Consequently, the maximum-likelihood fitting involved 5 free parameters in total: $\{w_{\text{early}}, w_{\text{late}}, \sigma_{\text{early}}, \sigma_{\text{late}}, \lambda\}$, where σ_{early} and σ_{late} are the standard deviations of Φ for early and late DM conditions, respectively.

To characterize the idiosyncratic patterns of stimulus-specific bias across participants, we defined the converging and diverging stimuli based on each individual's $\hat{\kappa}(\theta)$. We first identified the zero-crossing of $\hat{\kappa}(\theta)$ and estimated the local slopes. Then, for each participant, the nearby orientations (within $\pm 8^\circ$) were labeled as diverging (if the slopes were positive) and converging (if the slopes were negative) stimuli.

Quantifying the near-reference variability from behavior data

As a signature of the choice-induced bias, the marginal distribution of estimation error spreads more in near-reference trials than in far-reference trials (Figures S1D–S1K). We referred to this signature as the *near-reference variability* and characterized it by comparing the variability of the reference-conditioned estimation error distributions across the reference conditions. As a robust

measure of variability, we used the interquartile range (IQR), the difference between the first and third quartiles of the distribution. To capture the trend that the IQR increase as the reference nears the stimulus, we fitted a centered Gaussian density function, allowing the baseline, width, and amplitude parameters to vary.

To further validate our findings regarding the near-reference variability, we applied the same analysis procedure to publicly available datasets from previous studies (Figures S1A and S1B). In the work of Zamboni et al.¹⁰ and Fritsche and de Lange,⁸ the reference was used as a decision boundary as in our study, allowing us to assess whether the near-reference variability is a generalizable signature of the choice-induced bias and is well captured by our procedure. Additionally, to further confirm that the near-reference variability does not occur when a decision-making is not imposed as a task demand, even in the presence of a reference-like stimulus, we also applied the same analysis to another publicly available dataset from Rademaker et al.,⁷⁵ where the intervening orientation stimulus acts only as a distractor (Figure S1C). For the Rademaker et al.'s dataset, we included a relative stimulus range of $[-25^\circ, 25^\circ]$ in the analysis. Across all datasets, a reference was considered near if the relative stimulus orientation fell within $[-8^\circ, 8^\circ]$.

Quantifying the decision-consistent bias from behavior data

To characterize the decision-consistent bias, we computed the conditional mean of estimation errors ε given choice \hat{c} , denoting it as $b = (\mathbb{E}[\varepsilon|\hat{c} = cw] - \mathbb{E}[\varepsilon|\hat{c} = ccw])/2$. Previous studies^{2,9} showed that the decision-consistent bias is prominent only when the reference is near the stimulus orientation. Thus, we analyzed b only for the near-reference trials ($\hat{p} \in \{-4^\circ, 0^\circ, 4^\circ\}$) for further analyses.

Decomposition of the decision-consistent bias based on behavior data

We decomposed b into a component occurring before DM (pre-decision bias, b_{pre}) and the one after DM (post-decision bias, b_{post}):

$$b = b_{pre} + b_{post}. \quad (\text{Equation 4})$$

The decomposition was achieved in two steps, first quantifying the decision-consistent bias at the onset of DM from the discrimination data (b_{pre}) and then quantifying the additional decision-consistent bias accumulated after DM up until the moment of estimation (b_{post}). Conceptually, b_{pre} corresponds to the difference between the choice-conditioned means of memory states at the moment of discrimination time t_{dm} :

$$b_{pre} = (\mathbb{E}[m_{t_{dm}}|\hat{c} = cw] - \mathbb{E}[m_{t_{dm}}|\hat{c} = ccw]) / 2 \quad (\text{Equation 5})$$

Here, the underlying distribution of memory states $m_{t_{dm}}$ can be inferred from the parameters of the discrimination psychometric curve $\Psi(\theta, \rho, \vartheta)$ (defined in Equation 3), by applying the formalism offered by Signal Detection Theory.⁷⁶ Roughly put, this formalism relates the horizontal center and slope of Ψ to the mean and dispersion of the inferred distribution of memory states $m_{t_{dm}}$. Then the choice-conditioned means of this distribution, $\mathbb{E}[m_{t_{dm}}|\hat{c} = cw]$ and $\mathbb{E}[m_{t_{dm}}|\hat{c} = ccw]$, can readily be derived. We detailed this derivation in Method S1.1.

Next, having determined b_{pre} , we quantified b_{post} from the distribution of estimation errors ε , as follows. First, to enable the single-trial estimation, we first sign-flipped the estimation errors according to the choice direction, aligning their signs with b_{pre} , yielding sign-corrected errors $\varepsilon^* = \hat{c} \cdot \varepsilon$. We then subtracted b_{pre} from ε^* for each trial to compute the residuals, $\varepsilon^* - b_{pre}$. These residuals provide trial-to-trial estimates of how estimation errors are further deviated beyond b_{pre} . To quantify the decision-timing dependent changes in b_{post} , we performed a linear regression with condition indicators as regressors:

$$\mathbb{E}[\varepsilon^* - b_{pre}|\theta, \rho, \vartheta] = \beta_0 + \beta_1 \cdot \mathbf{1}_{late} + \beta_2 \cdot \hat{p} \quad (\text{Equation 6})$$

where β_0 corresponds to the b_{post}^{early} , β_1 to Δb_{post} , and β_2 was introduced as a nuisance parameter to capture the previously known attraction towards the reference⁷⁵.

For an additional comparison between the converging (conv) versus diverging (div) orientation conditions, we expanded Equation 6 to incorporate the converging-vs-diverging orientation factor, as follows:

$$\mathbb{E}[\varepsilon^* - b_{pre}|\theta, \rho, \vartheta] = \beta'_0 + \beta'_1 \cdot \mathbf{1}_{late} + \beta'_2 \cdot \mathbf{1}_{div} + \beta'_3 \cdot \mathbf{1}_{late,div} + \beta_4 \cdot \hat{p} \quad (\text{Equation 7})$$

where β'_0 corresponds to $b_{post}^{early,conv}$, β'_1 to Δb_{post}^{conv} , β'_2 to $b_{post}^{early,div} - b_{post}^{early,conv}$, β'_3 to $\Delta b_{post}^{div} - \Delta b_{post}^{conv}$, and β_4 to the reference attraction.

BOLD decoding of orientation memory states based on inverted encoding analysis

We decoded stimulus orientation from the population BOLD responses in the early visual cortex (V1, V2, and V3) using inverted encoding analysis.^{42,43} For each trial, the time courses of population BOLD responses \mathbf{X} ($N_{\text{voxels}} \times N_{\text{trials}}$) were modeled as a linear combination of channel responses \mathbf{Y} ($N_{\text{channels}} \times N_{\text{trials}}$) with weights \mathbf{W} ($N_{\text{voxels}} \times N_{\text{channels}}$), as follows:

$$\mathbf{X} = \mathbf{WY}, \quad (\text{Equation 8})$$

where each column of \mathbf{Y} corresponds to the vector of channel responses to stimulus orientation θ_j in a given trial j : $\mathbf{y}(\theta_j) = [y_{\psi_1}(\theta_j), \dots, y_{\psi_8}(\theta_j)]^\top$, where $y_{\psi}(\theta) = |\cos(\theta - \psi)|^8$ with channel centers ψ_1, \dots, ψ_8 tiling uniformly the orientation space $[0, \pi]$.

In the following steps, we carried out the decoding analysis using a leave-one-run-out cross-validation procedure. First, we designated one run as a held-out *validation* run and the remaining runs as *train* runs. Second, we constructed a matrix of population BOLD responses \mathbf{X}_T and a matrix of channel responses \mathbf{Y}_T from the train runs, along with a matrix of population BOLD responses \mathbf{X}_V from the

held-out validation run. Third, given \mathbf{X}_T and \mathbf{Y}_T , the weight matrix that yields the minimum squared errors $\widehat{\mathbf{W}}$ was determined as follows:

$$\widehat{\mathbf{W}} = \mathbf{X}_T \mathbf{Y}_T^T (\mathbf{Y}_T \mathbf{Y}_T^T)^{-1}. \quad (\text{Equation 9})$$

Fourth, we reconstructed the channel responses to the stimulus orientations presented in the validation run $\widehat{\mathbf{Y}}_V$ based on $\widehat{\mathbf{W}}$ and the matrix of population BOLD responses in the validation run \mathbf{X}_V , as follows:

$$\widehat{\mathbf{Y}}_V = (\widehat{\mathbf{W}}^T \widehat{\mathbf{W}})^{-1} \widehat{\mathbf{W}}^T \mathbf{X}_V. \quad (\text{Equation 10})$$

Fifth, to refine $\widehat{\mathbf{Y}}_V$ at a fine scale, we repeated the third and fourth steps while repositioning the centers of the eight channels, thereby defining $\widehat{\mathbf{Y}}_V$ for a total of 120 channel centers $\boldsymbol{\psi} = [0^\circ, 1.5^\circ, \dots, 178.5^\circ]^T$. We repeated the whole steps for each run, decision timing, and time point, resulting in the reconstruction of channel response vectors $\widehat{\mathbf{Y}}_V(t)$, with the columns corresponding to all the trials for a given time point t of interest 3–14 TRs.

From these reconstructed channel responses $\widehat{\mathbf{Y}}_V(t)$, we decoded the single-trial stimulus orientation by mapping the reconstructed channel response in each trial j and time point t ($\widehat{\mathbf{Y}}_j(t)$) to a point readout on the circular orientation space, as follows:

$$\widehat{\theta}_j^{bold}(t) = \text{atan} 2(\sin(2\boldsymbol{\psi})^T \widehat{\mathbf{Y}}_j(t), \cos(2\boldsymbol{\psi})^T \widehat{\mathbf{Y}}_j(t)) / 2. \quad (\text{Equation 11})$$

Estimating the time course of stimulus-specific bias in the memory states decoded from BOLD activity

The memory states decoded from BOLD activity exhibited the growth of the stimulus-specific bias over the delay. To track this growth, we estimated the amplitude of the bias at each time point based on the assumption that its across-stimulus profile is a scaled copy of the stimulus-specific bias estimated from the behavioral orientation estimates $\widehat{\kappa}(\theta)$, which was defined by the optimal parameters using Equation 1 for each participant. Accordingly, for each participant and each time point of BOLD measurement, we fitted the multiplicative weight of $\widehat{\kappa}(\theta)$ to the stimulus-specific errors of the memory states decoded from BOLD activity, $\widehat{\epsilon}^{bold}(t) = \widehat{\theta}^{bold}(t) - \theta$. For further analyses, we also considered the sign-corrected decoding errors $\widehat{\epsilon}^{*,bold}(t) = \widehat{c} \cdot \widehat{\epsilon}^{bold}(t)$, for which we flipped the signs of decoding errors $\widehat{\epsilon}^{bold}(t)$ for each trial according to the choice direction.

Decomposition of the decision-consistent bias based on BOLD activity

For the BOLD activity, we estimated the decision-consistent bias (b^{bold}) and its pre-decision (b_{pre}^{bold}) and post-decision (b_{post}^{bold}) components in the following steps. First, we modelled the time course of the latent, decision-consistent bias in memory states with a piecewise linear function $g(t)$ that incorporates the linear increase of the decision-consistent bias over time t along with the pulse-like shift due to the choice-induced bias during DM t_{dm} , as we assumed in our phenomenological models:

$$g(t) = (c_0 + c_1 \cdot t) \mathbf{1}_{t \leq t_{dm}} + (c_2^* + c_3 \cdot (t - t_{dm})) \mathbf{1}_{t > t_{dm}} \quad (\text{Equation 12})$$

where the first term on the right-hand side of captures the initial bias by c_0 , the linear increase over time by $c_1 \cdot t$ up to the time of DM, and the second term inherits the first term by including c_0 and $c_1 \cdot t_{dm}$ into $c_2^* = c_0 + c_2 + c_1 \cdot t_{dm}$ while capturing the pulse-like choice-induced bias by c_2 and the linear increase over time by $c_3 \cdot (t - t_{dm})$.

Second, we converted $g(t)$ to $\widehat{\epsilon}^{*,g(t)}$ using a transfer function ν , which convolves any given function with the canonical double-gamma hemodynamic response function⁷⁷ $h(t)$ while incorporating the input driven by the target stimulus θ and the reference stimulus ρ , which can be expressed as an argument of the complex number system as follows:

$$\widehat{\epsilon}^{*,g(t)} = \nu(g(t); \theta, \rho) = \arg(h(t) * (e^{2ig(t)} + \beta_\theta \cdot e^{2i\theta} \cdot \mathbf{1}_{t \in \mathcal{T}_\theta} + \beta_\rho \cdot e^{2i\rho} \cdot \mathbf{1}_{t \in \mathcal{T}_\rho})) / 2, \quad (\text{Equation 13})$$

where β_θ and β_ρ denote the beta weights of visual events⁴⁸ driven by the presentation of the stimulus θ and the reference ρ , while \mathcal{T}_θ and \mathcal{T}_ρ are the presentation time windows of θ and ρ .

Third, to estimate the influence of the stimulus and the reference (i.e., β_θ and β_ρ), we assumed that the impact of DM is negligible in $\widehat{\epsilon}^{*,bold}(t)$ (as defined at the end of the previous section) in the far-reference trials and defined its model correspondence by plugging the zero bias in the transfer function defined in Equation 13 instead of $g(t)$: $\widehat{\epsilon}^{*,0} = \nu(0; 0, \rho)$ (see Figure S4A). Then, we found the values of β_θ and β_ρ that minimize the L_2 difference between $\widehat{\epsilon}^{*,bold}(t)$ and $\widehat{\epsilon}^{*,0}$.

Fourth, having estimated β_θ and β_ρ , we then identified the parameters ($\widehat{c}_0, \widehat{c}_1, \widehat{c}_2, \widehat{c}_3$) of the time course of the latent, decision-consistent bias in memory states $g(t)$ that minimize the L_2 difference between $\widehat{\epsilon}^{*,bold}(t)$ in the near-reference trials and $\widehat{\epsilon}^{*,g}(t)$.

Lastly, we obtained the estimates of the decision-consistent bias (b^{bold}) and its pre-decision (b_{pre}^{bold}) and post-decision (b_{post}^{bold}) components, as follows:

$$\widehat{b}_{pre}^{bold} = \widehat{c}_0 + \widehat{c}_1 \cdot t_{dm}; \widehat{b}_{post}^{bold} = \widehat{c}_2 + \widehat{c}_3 \cdot (t_{em} - t_{dm}); \widehat{b}^{bold} = \widehat{b}_{pre}^{bold} + \widehat{b}_{post}^{bold} \quad (\text{Equation 14})$$

where t_{em} denotes the time of orientation estimation.

BOLD decoding of orientation memory states based on representational similarity analysis

As an alternative method for decoding orientation memory states from BOLD activity, we computed representational similarity matrices (RSMs) using the same population BOLD responses (\mathbf{X}) used in the inverted encoding analysis. Following previous studies,^{78,79} we averaged the across-voxel patterns of BOLD responses to each stimulus θ across trials (\mathbf{x}_θ) and then normalized \mathbf{x}_θ by subtracting the across-stimuli mean ($\bar{\mathbf{x}}_\theta$) from it. We defined the correlation matrix Ω ($N_{\text{stimulus}} \times N_{\text{stimulus}}$) by computing the Pearson correlation of $\mathbf{x}_\theta - \bar{\mathbf{x}}_\theta$ for each pair of orientation stimuli, while excluding identical cases that correspond to the diagonal elements. Such Ω was defined for each time point. We then estimated the orientation memory states ($\hat{\theta}_i^{sm}$) from Ω defined for each time point by taking the circular mean of each row, $\Omega_{i\cdot}$, corresponding to stimulus θ_i :

$$\hat{\theta}_i^{sm} = \text{atan} 2 \left(\sum_{j \neq i} \sin(2\theta_j) \Omega_{ij}, \sum_{j \neq i} \cos(2\theta_j) \Omega_{ij} \right) / 2 \quad (\text{Equation 15})$$

We quantified the “categorical” nature of orientation representation using a previously suggested category⁴⁰: whether a stimulus is on the clockwise or counter-clockwise side with respect to the vertical orientation. To model this categorical pattern, we used a “block” matrix $\mathbf{M}_{\text{block}}$ (each row is + 1 for within-category and − 1 for across-category) as well as a “cosine” matrix $\mathbf{M}_{\text{cosine}}$ (each row is the cosine value of the corresponding stimuli). To quantify the relative contribution of each pattern, we used the weight w_{cvx} between zero and one, such that the convex combination $\mathbf{M}_{\text{cvx}} = w_{\text{cvx}} \mathbf{M}_{\text{block}} + (1 - w_{\text{cvx}}) \mathbf{M}_{\text{cosine}}$ approximates Ω , indicating a higher categorical representation by a higher value of w_{cvx} . For comparison, we normalized Ω into the range $[-1, +1]$ and used the least squares method to find w_{cvx} for each time point.

Eye-tracking

To ensure participants’ eyes remained fixed on the central fixation marker throughout the experiment, we monitored their eye positions using an MR-compatible video-based eye tracker (EyeLink-1000, SR Research). The eye tracker was set up at a sampling rate of 500 Hz. For each participant, we recalibrated the eye tracker before each session using the built-in five-point routine (HV5). Eye-tracking data were corrupted or not recorded for five participants due to technical issues. Data was further excluded from analysis for the scan runs where experimenters noted calibration issues, which were attributable to eye occlusion by the head coil, unreliable tracking due to reflective sources like MRI goggles, or excessive blinking patterns.

Phenomenological models: diffusion-only and drift-diffusion models

Model description

We posited that the memory states in a single trial $m(t)$ undergo the following dynamics within the orientation space spanning $[0, \pi]$ with a periodic boundary:

$$m_t = m_0 + \int_0^t K(m_s) ds + \int_0^t D(m_s) dW_s + \alpha(m_{t_{\text{dm}}}, \rho) \cdot \mathbf{1}_{t \geq t_{\text{dm}}}, \quad (\text{Equation 16})$$

where $K(m_s)$ and $D(m_s)$ are the terms instantiating drift and diffusion dynamics, respectively, and W_s follows the Wiener process. We considered two classes of models, one with the diffusion term only (diffusion-only model) and the other with both drift and diffusion terms (diffusion-only model). The diffusion term $D(m)$, which is shared by both models, was set to w_D^2 . For the drift-diffusion model, the drift term $K(m)$ was defined in a stimulus-specific manner to instantiate the stimulus-specific drift, as follows:

$$K(m) = w_K \cdot \hat{\kappa}^\dagger(m), \quad (\text{Equation 17})$$

where w_K is the drift rate, and $\hat{\kappa}^\dagger$ is the normalized stimulus-specific bias, $\hat{\kappa}^\dagger = \hat{\kappa} / \max|\hat{\kappa}|$, with $\hat{\kappa}(\theta)$ defined by the optimal parameters found using (Equation 1) for each individual.

The last term of the right side of Equation 16 instantiates (i) the choice-induced bias by incurring an impulse-like shift in the memory trajectory in the choice-consistent direction and (ii) the reference-attraction bias at the moment of DM t_{dm} :

$$\alpha(m, \rho) = \begin{cases} w_\rho \cdot \tilde{\rho} + w_\alpha \cdot \hat{c}(m, \rho), & \tilde{\rho} \in \{-4^\circ, 0^\circ, 4^\circ\} \\ w_\rho \cdot \tilde{\rho}, & \tilde{\rho} \in \{-21^\circ, 21^\circ\} \end{cases}, \quad (\text{Equation 18})$$

where $\tilde{\rho}$ denotes the relative reference orientation; w_ρ is the strength of reference attraction, mimicking towards-distractor biases⁷⁵; w_α is the strength of choice-induced bias only present in the near reference conditions $\tilde{\rho} \in \{-4^\circ, 0^\circ, 4^\circ\}$, following the previous observations.^{2,9} The choice term \hat{c} was determined by the relative difference between the reference orientation ρ to the memory state at the moment of DM: $\hat{c} = \hat{c}(m, \rho) = \text{sign}(m - \rho)$.

To generate the discrimination and estimation reports, we used an instantaneous memory state at the corresponding moments of time: $m(t)$ at $t_{\text{dm}} = 6\text{s}$ and $t_{\text{dm}} = 12\text{s}$ to determine \hat{c} in the early and late DM conditions, respectively; $m(t)$ at $t_{\text{em}} = 18\text{s}$ to determine an estimation report.

Constraining the initial memory states based on the principle of efficient sensory encoding

We constrained the stimulus-specific distributions of initial memory states $p(m_0|\theta)$ based on the principle of efficient coding.²⁷ At the core of this principle is the encoding transformation function $F(\theta)$, which captures how encoding resources are allocated. $F(\theta)$ maps an orientation value θ in the stimulus space onto a measurement in the sensory space, in which the measurement is corrupted by the encoding noise. Therefore, using this framework, we first inferred $F(\theta)$ from data and then used it to derive $p(m_0|\theta)$.

We estimated the stimulus-to-sensory mapping F based on a previously derived relationship between the stimulus-specific bias $\kappa(\theta)$ and the derivative of $F(\theta)$ ²⁸:

$$F' \propto \left(\int \kappa d\theta \right)^{-1/2}. \quad (\text{Equation 19})$$

To compute the integral term in Equation 19, $\int \kappa d\theta$, we used previously estimated stimulus-specific bias $\hat{\kappa}(\theta)$ as an estimate for $\kappa(\theta)$. Given that $\hat{\kappa}(\theta) = \mathbf{v}(\theta)^\top \omega^*$, where $\mathbf{v}(\theta) = [1, v'_1(\theta), \dots, v'_{N_{\text{basis}}}(\theta)]^\top$ and ω^* is obtained from Equation 1, the integral becomes $\mathbf{V}(\theta)^\top \omega^*$, where $\mathbf{V}(\theta) = [v_1(\theta), \dots, v_{N_{\text{basis}}}(\theta)]^\top$ excluding the constant term. We then defined $\hat{\mathcal{K}}(\theta)$, our estimate of $\int \kappa d\theta$, with an additional adjustment by a shape parameter s in $[0, 1]$, which controls the extent to which the stimulus-specific bias constrains F' :

$$\hat{\mathcal{K}}(\theta) = (1 - s) \cdot \text{rescale}(\mathbf{V}(\theta)^\top \omega^*) + s, \quad (\text{Equation 20})$$

where $\text{rescale}(\cdot)$ denotes the min-max scaling between zero and one. As s approaches 0, the stimulus-to-sensory mapping becomes increasingly constrained by the integration of stimulus-specific bias function, and $s = 1$ corresponds to the case of uniform mapping. This estimator $\hat{\mathcal{K}}(\theta)$ allows us to specify F via Equation 19, by computing $F' \propto (\hat{\mathcal{K}}(\theta))^{-1/2}$ with the constraint $\int F'(\theta') d\theta' = \pi$.²⁸

Having specified the encoding transformation $F(\theta)$ and the sensory noise level w_E , we can determine the initial distribution of memory states $p(m_0|\theta)$ by modeling their initial states in the sensory space $F(m_0)$ as a von Mises distribution centered around $F(\theta)$ with dispersion proportional to w_E . The density function is computable using the change of variables:

$$p(m_0|\theta) = \frac{F'(m_0)}{2\pi I_0\left(\frac{1}{\sqrt{w_E}}\right)} \cdot \exp\left(\frac{1}{\sqrt{w_E}} \cdot \cos(2(F(m_0) - F(\theta)))\right) \quad (\text{Equation 21})$$

where F' is the derivative of F , and I_0 is the modified Bessel function of order 0. As such, we can fully constrain the stimulus-specific distribution of the initial memory states, $p(m_0|\theta)$, with the previously estimated $\hat{\kappa}$ and additional two parameters, s and w_E , which are fitted for both drift-diffusion and diffusion-only models (Figure S2).

Fitting the models to behavioral reports

To fit the models to the behavioral reports, we translated Equation 16 into the corresponding Fokker-Planck equation:

$$\frac{\partial}{\partial t} p(m, t) = - \frac{\partial}{\partial m} K^*(m, t) p(m, t) + \frac{w_D^2}{2} \frac{\partial^2}{\partial m^2} p(m, t), \quad (\text{Equation 22})$$

where $K^*(m, t) = K(m) + \alpha(m, \rho) \delta(t - t_{\text{dm}})$, where $\delta(\cdot)$ is the Dirac delta function. We numerically solved the equation by discretizing m with a unit $\Delta m = \pi/N_{\text{disc}}$, where $N_{\text{disc}} = 96$ (see Method S3.1 for detailed numerical procedure for model fitting). For each participant, we fit the models to the discrimination choices \hat{c}_j and the estimation reports $\hat{\theta}_j$ by finding the set of model parameters (8 parameters listed below) with the maximum joint likelihood $L(\text{parameters}|\text{data})$ given the experimental condition, which is specified by stimulus orientation θ_j , reference orientation ρ_j , and decision timing θ_j :

$$L(\text{parameters}|\text{data}) = \prod_{j=1}^{N_{\text{trial}}} p(\hat{c}_j, \hat{\theta}_j | \theta_j, \rho_j, \theta_j). \quad (\text{Equation 23})$$

We used the optimization routines provided by SciPy, with 20 iterations, while randomizing initial parameters by drawing from the constrained ranges of the model parameters (see Table S1). The parameters set free to be fitted were w_K (drift rate) (set to 0 for the diffusion-only model), w_D (diffusion rate), w_α (choice-induced bias), w_ρ (reference bias weight), w_E (encoding variability), w_P (production variability), w_λ (decision-making lapse), and s (encoding function shape). To compute cross-validated log likelihoods, we ran 10 independent runs of 5-fold cross-validation of log likelihoods (each with 20 iterations) by separating the data used for fitting the models including the estimation of $\hat{\kappa}$.

Evaluating the correspondence between the drift-diffusion model and BOLD activity

We validated the drift-diffusion model's prediction of the memory state dynamics by evaluating how closely it follows the trajectories of the memory states decoded from BOLD activity. For this evaluation, we translated the model prediction of memory states m_t , as defined in Equation 16, into its equivalent in BOLD signal using the transfer function $\nu(m_t; \theta, \rho)$, as defined in Equation 13. Then, for each trial j in the near-reference conditions, we evaluated the correspondence between the model prediction $\nu(m_j; \theta_j, \rho_j)$ and the

memory states decoded from BOLD activity $\hat{\theta}_j^{bold}$ by quantifying the average cosine distance between them with a correspondence score S_j , as follows:

$$S_j = \mathbb{E} \left[\cos \left(2 \left(\hat{\theta}_j^{bold} - \nu(m_j; \theta_j, \rho_j) \right) \right) \right]. \quad (\text{Equation 24})$$

Recurrent neural network model

RNN dynamics

The following equations describe the dynamics of RNN:

$$\tau \frac{d\mathbf{r}}{dt} = -\mathbf{r} + f(\mathbf{J}\mathbf{r} + \mathbf{J}_\theta \mathbf{u}_\theta + \mathbf{J}_\rho \mathbf{u}_\rho + \boldsymbol{\eta}); \quad (\text{Equation 25})$$

$$\tau' \frac{d\boldsymbol{\eta}}{dt} = -\boldsymbol{\eta} + \sqrt{2\tau'} \boldsymbol{\xi}, \quad (\text{Equation 26})$$

where \mathbf{r} is the N_{rec} -dimensional ($N_{\text{rec}} = 96$) unit activity with a time constant $\tau = 100$ ms, and $\boldsymbol{\eta}$ is the stochastic noise with a time constant $\tau' = 200$ ms, modeled as the Ornstein–Uhlenbeck process⁸⁰; \mathbf{u}_θ and \mathbf{u}_ρ are the N_{in} -dimensional ($N_{\text{in}} = 24$) stimulus and reference inputs, respectively, whose units are orientation tuned, modeled as a von Mises distribution; \mathbf{J} , \mathbf{J}_θ and \mathbf{J}_ρ are the weights of the recurrent, stimulus, and reference inputs, respectively; $f = 1/(1 + \exp(-x))$ is the sigmoid activation function; $\boldsymbol{\xi}$ is the independent Gaussian noise with a standard deviation of 0.05. The values for \mathbf{r} and $\boldsymbol{\eta}$ were initialized at 0s. We approximated the equations above using the forward Euler approximation with a discretization time step $\Delta t = 20$ ms.

Considering that the stimulus and reference inputs occupied different parts of the visual field in the task paradigm, \mathbf{u}_θ and \mathbf{u}_ρ were projected onto two separate 48-dimensional populations of the recurrent units \mathbf{r} , namely \mathbf{r}_θ and \mathbf{r}_ρ . \mathbf{u}_θ is centered at veridical orientation θ , given as

$$(\mathbf{u}_\theta)_i = \gamma_\theta \cdot \exp(\kappa_\theta (\cos(2(\theta - \theta_i)) - 1)), \quad (\text{Equation 27})$$

where θ_i is the preferred orientation of the unit i ; γ_θ is the strength of the stimulus input, fixed at 1; κ_θ is the concentration parameter, fixed at 5. \mathbf{u}_ρ during the discrimination epoch was determined by a one-hot vector:

$$(\mathbf{u}_\rho)_i = \gamma_\rho \cdot \delta_{\rho, \rho_i}, \quad (\text{Equation 28})$$

where ρ is the reference orientation, and ρ_i is the preferred orientation of the unit i ; γ_ρ is the strength of the reference input, fixed at 2, which is higher than the one for \mathbf{u}_θ considering the higher level of certainty. We used $N_\theta = 24$ ranging from 0° to 172.5° with 7.5° increments. The reference input was constrained to $|\rho - \theta| \leq 30^\circ$ with 7.5° steps, resulting in 9 possible relative references. During the “train episode,” we excluded $\rho = \theta$ to facilitate training but included it during the “generalization episode” (see the next section for the definition of the “train episode” and “generalization episode”).

Discrimination and estimation outputs, \mathbf{z}^{dm} and \mathbf{z}^{em} , were

$$\mathbf{z}^{dm} = \mathbf{J}_{dm} \mathbf{r}, \mathbf{z}^{em} = \mathbf{J}_{em} \mathbf{r}, \quad (\text{Equation 29})$$

where $\mathbf{z}^{dm} = (z_1^{dm}, z_2^{dm})$ for CW and CCW choices, respectively, and \mathbf{z}^{em} , a 24-dimensional ‘labeled line’ response vector, consists of equally discretized points within $[0, \pi]$. Input and output weights, \mathbf{J}_θ , \mathbf{J}_ρ , \mathbf{J}_{em} , and \mathbf{J}_{dm} , were fixed, defined as

$$\mathbf{J}_\theta(j, k) = \mathbf{J}_\rho(j, k) = \gamma_{in} \cdot \cos(2\pi(j - k)/N_{in}); \quad (\text{Equation 30})$$

$$\mathbf{J}_{em}(j, k) = \gamma_{em} \cdot \cos(2\pi(j - k)/24); \quad (\text{Equation 31})$$

$$\mathbf{J}_{dm}(j, k) = \mathbf{1}_{j \in I_{cw}} \cdot \mathbf{1}_{k=1} + \mathbf{1}_{j \in I_{ccw}} \cdot \mathbf{1}_{k=2}, \quad (\text{Equation 32})$$

where $\gamma_{in} = 1$; $\gamma_{em} = 0.4$; a balanced voting for the two choices, with $I_{cw} = \{j + 1 : (j \bmod 48) < 24\}$ and $I_{ccw} = \{j + 1 : (j \bmod 48) \geq 24\}$.

RNN training

We trained RNNs using a paradigm equivalent in structure to the one used for human participants. RNNs were first trained on a short task timescale (“train episode”) and then generalized to an extended timescale (“generalization episode”). As we confine RNNs as a proof of principle, the time scales of dynamics were not directly aligned with the human experiment. Each trial in the train episode had the following structure: an initial fixation epoch with no inputs (0.1s) was followed by stimulus presentation (0.6s), first delay (0.3s), DM (0.6s), second delay (0.3s), and estimation report (0.1s) epochs. In the generalization episode, respecting the human task structure, the first and second delays were extended to 1.8 s and 4.2 s for the early DM condition, and 4.2s and 1.8 s for the late DM condition, while the lengths of the other epochs remained the same.

For the supervised learning, we defined desired outputs \mathbf{q}^{dm} and \mathbf{q}^{em} as

$$\mathbf{q}^{dm} = [\mathbf{1}_{\theta > \rho}, \mathbf{1}_{\theta < \rho}]^T; \quad (\text{Equation 33})$$

$$(\mathbf{q}^{em})_j = \exp(\kappa_\theta (\cos(2(\theta - \theta_j)) - 1)), \quad (\text{Equation 34})$$

where \mathbf{q}^{dm} and \mathbf{q}^{em} were 2-dimensional and 24-dimensional vectors, respectively. We trained the recurrent weight \mathbf{J} while maintaining other weights fixed. Before training, \mathbf{J} was initialized as zero. The joint loss \mathcal{L} was $\mathcal{L}_{dm} + \mathcal{L}_{em}$, where both \mathcal{L}_{dm} and \mathcal{L}_{em} are the time-averaging cross entropies between the network output and the desired output, given as:

$$\mathcal{L}_{dm} = \left\langle \sum_{j \in \{\text{cw}, \text{ccw}\}} M^{dm}(t) \cdot \mathbf{z}_j^{dm}(t) \cdot \log\left(1 / \mathbf{q}_j^{dm}(t)\right) \right\rangle_t; \quad (\text{Equation 35})$$

$$\mathcal{L}_{em} = \left\langle \sum_{j \in \{1, \dots, 24\}} M^{em}(t) \cdot \mathbf{z}_j^{em}(t) \cdot \log\left(1 / \mathbf{q}_j^{em}(t)\right) \right\rangle_t, \quad (\text{Equation 36})$$

where $M^{dm}(t)$ is a binary mask, non-zero only during the discrimination epoch, while $M^{em}(t)$ is non-zero except for the initial fixation epoch. The loss was minimized using backpropagation in PyTorch with the Adam optimizer (with a learning rate of 0.02). We undertook 300 iterations per network training, generating 128 trials per iteration. In each of those trials, the stimulus and relative reference orientations were determined randomly.

To dissociate the effects of drift dynamics and the choice-induced bias, we independently trained 50 “homogeneous” RNNs, along with the original “heterogeneous” RNNs. For the heterogeneous RNNs (Figures 6 and 8), to approximate the orientation-specific variability that reflects the efficient sensory encoding principle, we added Gaussian noise to the orientation input θ , allowing the centers of \mathbf{u}_θ for a given stimulus orientation θ_0 to vary as follows:

$$p(\theta|\theta_0) \sim \mathcal{N}(\theta_0, \gamma_D \cdot |\sin(2\theta_0)|^2), \quad (\text{Equation 37})$$

where the spread term γ_D was set at 10^2 . In contrast, no such noises were added for the homogeneous RNNs ($\theta|\theta_0 = \theta_0$; Figure 7). All training details, except for stimulus input-level encoding variability, were identical for the homogeneous and heterogeneous RNNs.

To inspect the effect of feedback connections on the choice-induced bias, we independently trained 50 feedback-connection-ablated RNNs by zeroing the connectivity from the units receiving \mathbf{u}_ρ to those receiving \mathbf{r}_θ . To examine the effect of fine-tuning the readout connection after training the feedback-ablated RNNs, we further trained the readout connection independently (mapping from recurrent activities \mathbf{r} to both discrimination and estimation outputs \mathbf{z}^{dm} and \mathbf{z}^{em}).

RNN analysis

From the output vectors \mathbf{z}^{dm} and \mathbf{z}^{em} of the 50 independently trained RNNs, we determined their discrimination choice \hat{c}^{mn} and estimation report $\hat{\theta}^{mn}$, as follows:

$$\hat{\theta}^{mn} = \text{atan} 2 \left(\sum_{j=1}^{24} \mathbf{z}_j^{em} \sin 2\theta_j, \sum_{j=1}^{24} \mathbf{z}_j^{em} \cos 2\theta_j \right) / 2; \quad (\text{Equation 38})$$

$$\hat{c}^{mn} = \text{sign}(\mathbf{z}_1^{dm} - \mathbf{z}_2^{dm}). \quad (\text{Equation 39})$$

We then conducted the same analyses on \hat{c}^{mn} and $\hat{\theta}^{mn}$, as we did on human discrimination choices \hat{c} and estimation reports $\hat{\theta}$, to assess the stimulus-specific and decision-consistent biases.

To depict RNNs’ population dynamics during the discrimination epoch in a low-dimensional state space, we applied PCA on average RNNs, taking the mean of individually trained \mathbf{J} for the homogeneous and heterogeneous RNNs separately. We generated trials from both types of RNNs without network noise (i.e., $\xi = 0$ in Equation 26) for this state-space analysis. We analyzed the dynamics and connectivity patterns by further separating both \mathbf{r}_θ and \mathbf{r}_ρ into the CW- and CCW-projecting populations based on \mathbf{J}_{dm} (e.g., \mathbf{r}_ρ^{cw} denotes the reference-receiving and CW-projecting population). For illustration, we used a near reference, $\theta - \rho = 7.5^\circ$, with winning/losing populations as $\mathbf{r}^{cw}/\mathbf{r}^{ccw}$.

To inspect the dynamics of the homogeneous and heterogeneous RNNs, we projected the activities of mean homogeneous and heterogeneous RNNs onto the principal components from mean homogeneous RNNs, assuming slow drift dynamics in heterogeneous RNNs. We first stacked the activities of $\bar{\mathbf{r}}_\theta = (\mathbf{r}_\theta^{cw} + \mathbf{r}_\theta^{ccw})/2$ of the mean homogeneous RNNs for each condition ($N_\theta = 24$ stimulus orientations and $N_\rho = 9$ relative reference orientations) to form a column-mean-centered data matrix \mathcal{D} , that is, $(N_\theta \cdot N_\rho \cdot T) \times 24$, where T is the total time steps. In the text, \mathbf{r}_θ is used in the place of $\bar{\mathbf{r}}_\theta$ for brevity. The PC projection matrix $\mathbf{V}_\mathcal{D}$ was computed from the singular value decomposition of \mathcal{D} as $\mathcal{D} = \mathbf{U}_\mathcal{D} \mathbf{S}_\mathcal{D} \mathbf{V}_\mathcal{D}^T$. We used the first two PCs to project population activities of \mathbf{r}_ρ^{cw} , \mathbf{r}_ρ^{ccw} , and $\bar{\mathbf{r}}_\theta$. The first two PCs explained more than 92% of the total variance. For intuitive presentation, we sign-flipped and rotated the projection axes to align different stimulus conditions clockwise, with 0° stimulus points upright in \mathbf{r}_θ space, resulting in oblique stimuli ($\pm 45^\circ$) along the x-axis and cardinal stimuli ($0^\circ, 90^\circ$) along the y-axis.

QUANTIFICATION AND STATISTICAL PROCEDURES

For the quantitative evaluations of phenomenological models (Figures 4D, 4E, and 5F), we simulated the trajectories of memory states based on 10,000 Monte-Carlo iterations. For the statistical analyses of population-level differences in the decision-consistent biases around converging and diverging fixed points in the human behavior and drift-diffusion model (Figures 4E and 4F), the BOLD signals (Figure 5E), and RNN models (Figure 6H), we ran bootstrap-based permutation test using 10,000 random iterations.