

1 Sensory readout accounts for adaptation

Timothy C. Sheehan¹, John T. Serences^{1,2,3}

¹Neurosciences Graduate Program, ²Department of Psychology, and ³Kavli Institute for Brain and Mind, University of California, San Diego, California 92093

2 Abstract

3 Sensory responses and behavior are strongly shaped by stimulus history. For instance, perceptual
4 reports are sometimes biased towards previously viewed stimuli (*serial dependence*). Previous
5 behavioral studies suggest that serial dependence is implemented via modulations in visual cortex, but
6 neural evidence is lacking. We recorded fMRI responses while human participants performed a delayed
7 orientation discrimination task. While behavioral reports were *attracted* to the previous stimulus,
8 response patterns in sensory areas were *repelled*. We reconciled these opposing biases using a model
9 where both sensory encoding and readout are shaped by stimulus history. Neural adaptation reduces
10 redundancy at encoding and leads to the repulsive biases that we observed in visual cortex. Serial
11 dependence is not implemented in visual cortex but rather by readout mechanisms that account for
12 adaptation during encoding. The model suggests the visual system improves efficiency via adaptation
13 while still optimizing behavioral readout based on the temporal structure of natural stimuli.

14 Keywords: neural adaptation, sensory decoding, Bayesian inference, serial dependence

15 Introduction

16 Natural stimuli are known to demonstrate statistical dependencies across both space and time,
17 such as a prevalence of vertical and horizontal (cardinal) orientations and a higher probability of small
18 orientation changes in sequential stimuli¹⁻³. These regularities can be leveraged to improve the
19 efficiency and accuracy of visual information processing. We use the term “encoding” to refer to the
20 initial conversion of external sensory information into neural activity patterns and the term “readout” to
21 refer to the readout of these encoded signals to shape behavior. At encoding, regularities yield
22 attenuated neural responses to frequently occurring stimuli (“adaptation”), reducing metabolic cost and
23 redundancy in neural codes. At readout, regularities support the formation of Bayesian priors that can
24 be used to bias perception in favor of higher probability stimuli. On their own, both adaptation and
25 Bayesian readout can explain a variety of behavioral phenomena such as improved precision around the
26 cardinal axes (oblique effect) and why we remember objects as being closer to the average exemplar
27 (contraction bias)⁴⁻⁶. While the effects of sensory history on sensory coding and behavior have been
28 studied extensively, it is unclear how changes at encoding shape readout and behavior.

29 Adaptation increases coding efficiency by modulating sensory tuning properties as a function of
30 the recent past. For instance, reducing the gain of neurons tuned to a recently seen adapting stimulus
31 reduces the temporal autocorrelation of activity when similar stimuli are presented sequentially. In turn,
32 reducing these autocorrelations improves the overall efficiency of sensory codes: fewer spikes are
33 dedicated to encoding redundant stimuli, and the presence of a novel stimulus can be more easily
34 detected as it will be accompanied by a sudden increase in activity^{2,7-15}. Importantly, adapted
35 representations early in the processing stream (e.g. in LGN) are inherited by later visual areas^{11,16,17}.

36 Although adaptation increases coding efficiency, it comes at a cost to perceptual fidelity as
37 adaptation can lead to repulsion away from the adapting stimulus for features such as orientation and
38 motion direction^{18–20}. For instance, after continuously viewing and adapting to motion in one direction,
39 stationary objects will appear to be moving in the opposite direction (i.e., current perceptual
40 representations are *repelled* away from recent percepts). However, this potentially deleterious
41 aftereffect is accompanied by better discriminability around the adapting stimulus, which may be more
42 important than absolute fidelity from a fitness perspective^{14,21–24}.

43 In contrast to the repulsive biases associated with neural adaptation, perception is sometimes
44 attracted to recently attended items in conditions where weak stimuli are attended – a phenomenon
45 termed “serial dependence”. As serial dependence can impact immediate perceptual reports and the
46 relative perception of simultaneously presented items, some have suggested that it reflects modulations
47 in early stages of sensory processing^{25–27}. In line with this idea, one fMRI study demonstrated that early
48 sensory biases match ‘attractive’ behavioral reports²⁸. However, consecutive stimuli were always either
49 the same or orthogonal orientations, conditions where serial dependence effects on behavior are
50 negligible^{26,29,30}. Thus, without sampling the entire stimulus feature space, it is unclear how to integrate
51 this finding with related empirical and theoretical work on serial dependence. Counter to studies
52 suggesting a sensory locus of serial dependence, other behavioral results have found that serial
53 dependence does not occur immediately after encoding but instead emerges only, and increases with, a
54 working memory maintenance period^{31–33}. This observation suggests that serial dependence is not the
55 product of early sensory coding³² and instead might be implemented by a later readout or memory
56 maintenance circuit^{34,35}. There is evidence that such a readout mechanism is Bayesian, as the influence
57 of the “prior” (the previous stimulus) is larger when sensory representations are less precise due to
58 either external or internal noise^{27,35}. Thus, the collective evidence is mixed, with some studies pointing
59 towards an early sensory locus and others to later stages of readout and memory storage.

60 This lack of consensus suggests that assessing interactions between sensory and readout stages
61 of processing may be key to better understanding the impact of stimulus history on perception. For
62 example, previous work suggests that readout does not account for neural adaptation that happens
63 during encoding, as adaptive repulsive biases cascade across layers of the visual processing hierarchy
64 and penetrate behavioral reports^{11,17,36}. These studies, however, did not consider paradigms where the
65 adapting stimulus was behaviorally relevant. Attending to relevant stimuli may shape how readout
66 stages account for the current state of adaptation, possibly inducing attractive serial dependence. To
67 assess this possibility, we utilized multivariate fMRI decoding techniques to characterize how
68 representations in early visual areas change as a function of stimulus history during a delayed
69 orientation discrimination task (Figure 1A). We replicated classic “serial dependence” findings where
70 behavioral reports were attracted to the orientation of the previous stimulus. We found that this
71 attractive behavioral bias was not accompanied by attractive biases in visual cortex, as predicted by
72 early sensory models of serial dependence. Rather, we observed *repulsive biases* in early visual cortex,
73 consistent with adaptation. To explain these results, we examined several possible read-out
74 mechanisms and found that only decoding schemes that account for adaptation can explain attractive
75 serial dependence in behavior. More generally, these results explain how the visual system can reduce
76 energy usage without sacrificing precision by optimizing encoding and behavioral readout relative to the
77 temporal structure of natural environments.

78 Results

79 Behavior

80 To probe the behavioral effects of serial dependence, we designed a delayed discrimination task
 81 where participants judged whether a bar was tilted CW or CCW relative to the orientation of a
 82 remembered grating (Figure 1A). We first report the results from a behavior-only study (n=47) followed
 83 by an analysis of neural activity for a cohort completing the same task in the scanner (n=6). Task
 84 difficulty was adjusted for each participant by changing the magnitude of the probe offset ($\delta\theta$) from the
 85 remembered grating and was titrated to achieve a mean accuracy of $\sim 70\%$ (accuracy $69.8 \pm 0.82\%$, $\delta\theta$:
 86 $4.95 \pm 0.27^\circ$; all reported values mean $\pm 1\text{SEM}$ unless otherwise noted). Fixing subjects at this
 87 intermediate accuracy level helped to avoid floor/ceiling effects and improved our sensitivity to detect
 88 perceptual biases while keeping participants motivated.

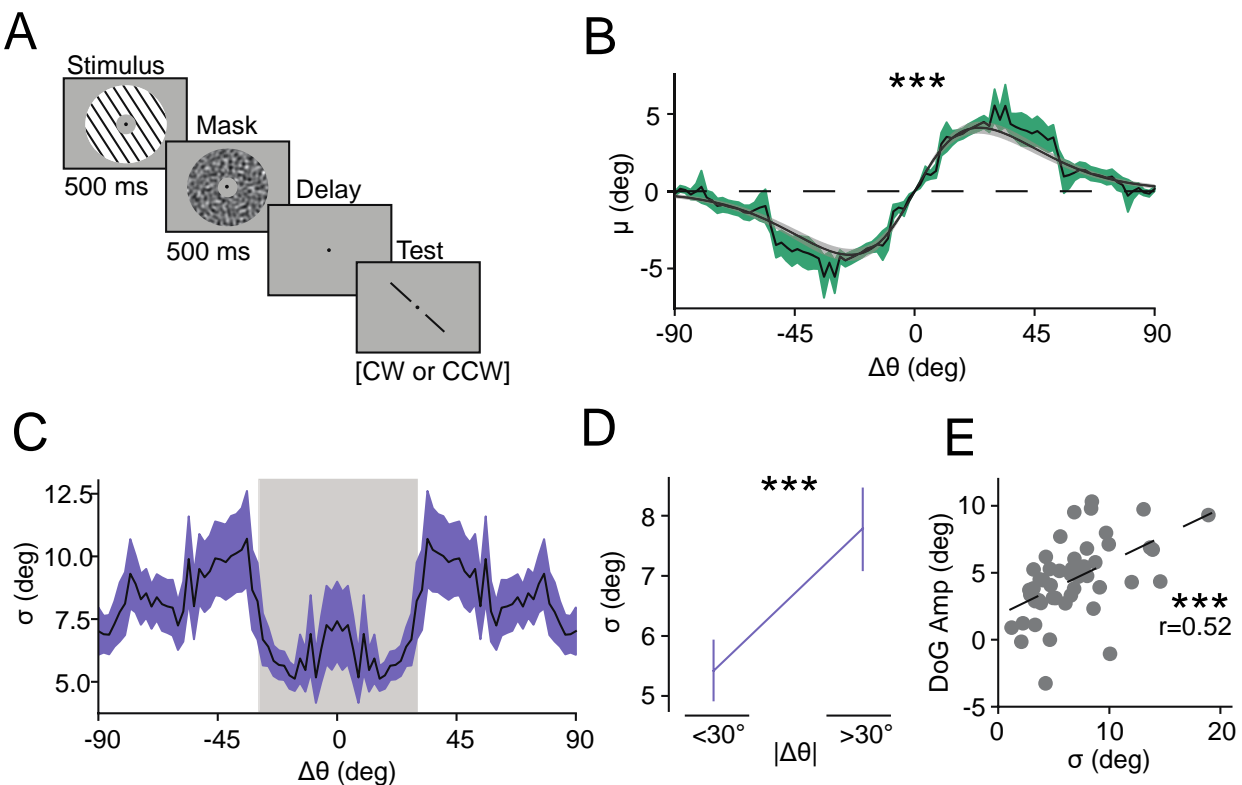


Figure 1 Behavior. **A:** Task Schematic. An orientated stimulus is followed by a probe bar that is rotated $<15^\circ$ from the stimulus. Participants judged whether the bar was CW or CCW relative to the stimulus in a binary discrimination task. **B:** Behavioral bias, green: average model-estimated bias as a function of $\Delta\theta = \theta_{n-1} - \theta_n$ ($\pm 1\text{SEM}$ across participants); gray: average DoG fit to raw participant responses sorted by $\Delta\theta$ ($\pm 1\text{SEM}$ across participants). **C:** Behavioral σ as a function of $\Delta\theta$, shaded region corresponds to $|\Delta\theta| < 30^\circ$. **D:** Behavioral variance is significantly less for $|\Delta\theta| < 30^\circ$. **E:** Bias is correlated with variance across participants. ***, $p < .001$

89 To quantify the pattern of behavioral responses, we modelled the data as the product of a noisy
 90 encoding process described by a Gaussian function centered on the presented orientation with standard
 91 deviation σ and bias μ . Optimal values for σ and μ were found by maximizing the likelihood of responses
 92 for probes of varying rotational offsets from the remembered stimulus, thus converting pooled binary
 93 responses into variance and bias measured in degrees (see [Response Bias](#), Figure S1). This allowed us to
 94 measure precision for individual participants and also allowed us to measure how responses were

95 biased as a function of the orientation difference between the remembered gratings on consecutive
96 trials $\Delta\theta = \theta_{n-1} - \theta_n$, an assay of serial dependence. To increase power and remove systematic response
97 biases, we ‘folded’ responses relative to $\Delta\theta$ such that all analysis have either rotational (Figure 1B) or
98 horizontal (Figure 1C) symmetry³⁴.

99 Responses were clearly biased towards the previous stimulus (Figure 1B, green curve), which
100 we quantified by fitting a Derivative-of-Gaussian (DoG) function to the raw response data for each
101 participant (gray curve; amplitude: $4.53^\circ \pm 0.42^\circ$, $t(46) = 7.8$, $p = 5.9 \times 10^{-10}$, one sample t-test; full width at
102 half max (FWHM): $42.9^\circ \pm 1.8^\circ$, see [Serial Dependence](#)). The magnitude and shape of serial dependence is
103 consistent with previous reports^{26,29}. We next examined how response precision (σ) varied as a function
104 of $\Delta\theta$ and found that responses were more precise around small changes (Figure 1C), again consistent
105 with previous reports³⁷. We quantified this difference in precision by splitting trials into ‘close’ and ‘far’
106 bins (greater than or less than 30° separation) and confirmed that responses following ‘close’ stimuli
107 were more precise ($t(46) = -3.72$, $p = 0.0003$, paired 1-tailed t-test, Figure 1D, see [Response Precision](#)).
108 Note that the choice of 30° was arbitrary and all threshold values between 20° and 40° yielded a
109 significant ($p < .05$) result.

110 Previous work has shown that serial dependence is greater when stimulus contrast is lower²⁷
111 and when internal representations of orientation are weaker due to stimulus independent fluctuations
112 in encoding fidelity³⁵. We tested a Bayesian interpretation of these findings by asking whether less
113 precise individuals are more reliant on prior expectations and therefore more biased. Indeed, we found
114 a positive correlation between DoG amplitude and σ , (Figure 1E, $r(45) = 0.52$, $p = .0001$, 1-tailed Pearson’s
115 correlation).

116 Finally, we confirmed that our behavioral results were not driven by an artifact of our modeling
117 procedure by replicating the relationships using raw response probabilities (Figure S2). We also
118 confirmed these effects were not driven by a subset of low/high performing individuals or trial counts
119 (Figure S3) and were not due to inhomogeneities in the stimulus sequences used for some participants
120 (Figure S4).

121 [Stimulus history effects in visual cortex](#)

122 To examine the influence of stimulus history on orientation-selective response patterns in early
123 visual cortex, six participants completed between 748 and 884 trials (mean 838.7) of the task in the fMRI
124 scanner over the course of four, two-hour sessions (average accuracy of $67.7\% \pm 0.4\%$ with an average
125 probe offset, $\delta\theta$, of 3.65°). As with the behavior-only cohort, these participants showed strong attractive
126 serial dependence (Figure 2A, green) that was significantly greater than 0 when parameterized with a
127 DoG function (amplitude = $3.50^\circ \pm 0.27^\circ$, $t(5) = 11.93$, $p = .00004$; FWHM = $35.9^\circ \pm 2.34^\circ$, Figure 2A black dotted
128 line). This bias was not significantly modulated by inter-trial interval, delay period, or an interaction
129 between the two factors (all p-values > 0.5 , mixed linear model grouping by participant). Similar to the
130 behavioral cohort, we found that variance was generally lower around small values of $\Delta\theta$. We quantified
131 variance in the same manner as the behavioral cohort (flipping responses to match biases and down-
132 sampling the larger group) and found that responses were more precise following close ($< 30^\circ$) relative to
133 far stimuli ($> 30^\circ$, $t(5) = -9.96$, $p = 0.00009$, 1-tailed paired t-test, Figure 2B). This pattern was significant
134 ($p < 0.05$) for thresholds between 20° and 40° , and these findings were generally replicated when

135 analyzed using a model-free approach as described for the behavior-only cohort (Figure S2D-F, except
136 that the threshold analysis no longer reached significance).

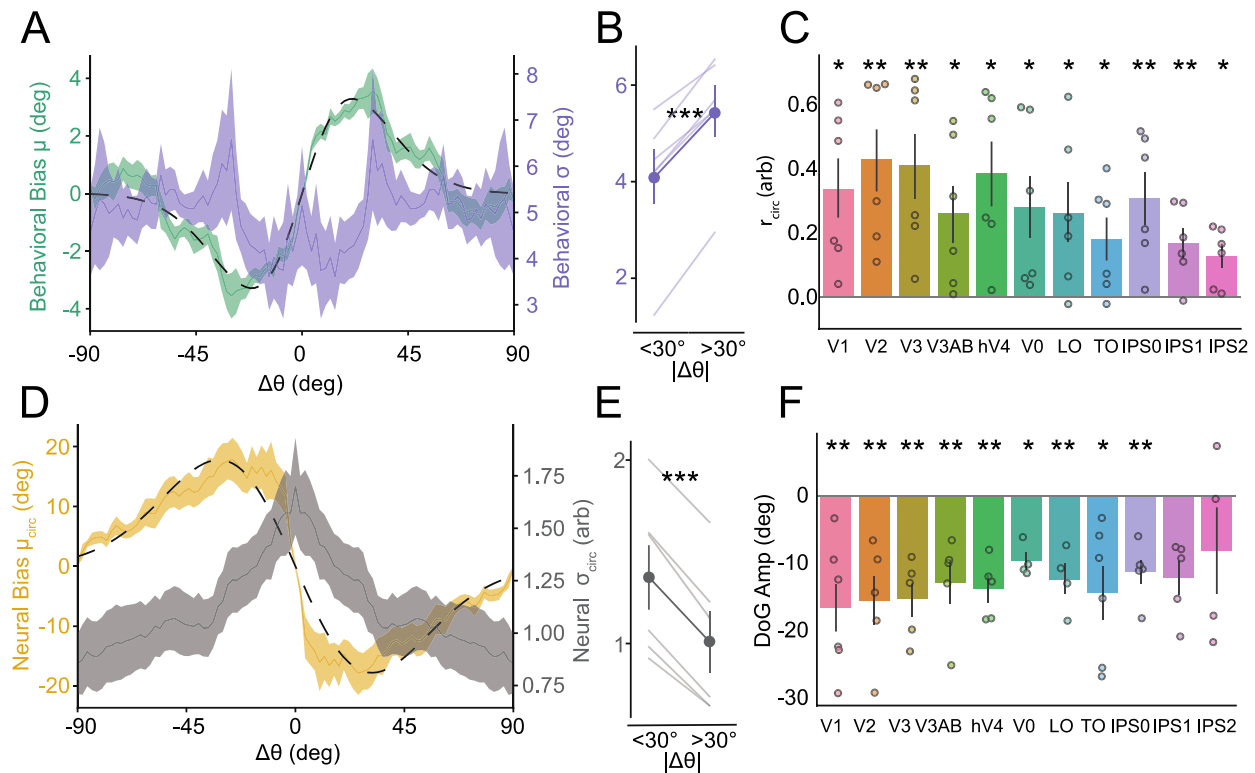


Figure 2 Behavioral and Neural Bias. **A:** Left-axis, Behavioral serial dependence. Shaded green: average model-estimated bias as a function of $\Delta\theta$ (\pm SEM across participants); dotted black line: average DoG fit to raw participant responses sorted by $\Delta\theta$. Right-axis, variance. Purple shaded line: model-estimated variance as a function of $\Delta\theta$ (\pm SEM across participants). **B:** Behavioral σ is significantly less for $|\Delta\theta| < 30^\circ$. **C:** Decoded orientation was significantly greater than chance when indexed with circular correlation for all ROIs examined. Error bars indicate \pm SEM across participants. Dots show data from individual participants. **D:** Left-axis, decoding bias. Shaded yellow line: decoded bias (μ_{circ} of decoding errors) sorted by $\Delta\theta$ (\pm SEM across participants); dotted black line: average DoG fit to raw decoding errors sorted by $\Delta\theta$. Right-axis, decoded σ_{circ} . Shaded gray line: average decoding variance (σ_{circ}) as a function of $\Delta\theta$ (\pm SEM across participants). Note that σ_{circ} can range from $[0, \infty]$ and has no units. **E:** Decoded variance is significantly less for $|\Delta\theta| < 30^\circ$. **F:** Decoded errors are significantly repulsive when parameterized with a DoG for most ROIs. *, $p < .05$; **, $p < .01$; ***, $p < .001$.

137 To characterize activity in early visual areas, independent retinotopic mapping runs were
138 completed by each subject to identify regions of interest (ROIs) consisting of: V1, V2, V3, V3AB,
139 Ventral-, Temporal- and Lateral-Occipital Areas (VO, TO, and LO), and intraparietal sulcus areas IPS0-2. In
140 addition, a separate localizer task was used to sub-select the voxels that were most selective for the
141 spatial position and orientation of the stimuli used in our task (see [Voxel Selection](#)).

142 To examine how visual representations are affected by stimulus history, we trained a decoder
143 on the orientation of the sample stimulus on each trial based on BOLD activation patterns in each ROI.
144 We used a decoder that accounts for stimulus related noise correlations (see [Orientation Decoding](#),^{35,38})
145 using a leave-one-run-out cross-validation across sets of 68 consecutive trials (4 blocks of 17 trials) that
146 had orientations pseudo randomly distributed across all 180° of orientation space. We first quantified

147 single-trial decoding performance using circular correlation (r_{circ}) between the decoder-estimated
148 orientations and the actual presented orientations and found that all ROIs had significant orientation
149 information (Figure 2C). For this and all of the remaining main figures, we used the average of four TRs
150 (spanning 4.8-8.0s) following stimulus presentation to avoid the influence of the probe stimulus (which
151 came up $\geq 6s$ into the trial and thus should not influence responses in the 4.8-8.0s window after
152 accounting for hemodynamic delay). That said, our ability to decode orientation was not specific to the
153 exact TRs selected, or whether the decoder was trained on the task data or an independent localizer
154 task (Figures S5-S6).

155 The high SNR of the BOLD decoder allowed us to examine residual errors on individual trials.
156 When measuring the bias of these decoding errors (μ_{circ}) as a function of stimulus history ($\Delta\theta$), we
157 unexpectedly observed a strong *repulsive* bias reflecting neural adaptation (V3, Figure 2G yellow, see
158 [Neural Bias](#)). This bias was significant when quantified with a DoG (amplitude= $-19.76^\circ \pm 3.06^\circ$, $t(5)=-5.90$,
159 $p=.0010$; FWHM= $47.7^\circ \pm 1.71^\circ$, Figure 2G black dotted-line). All ROIs except IPS0 and IPS1 had a
160 significantly negative amplitude ($p<.05$) and the average DoG amplitude across ROIs was also significant
161 ($t(10)=-7.65$, $p=.00001$, Figure 2F). This repulsive pattern suggests that serial dependence is not a direct
162 result of biases in early sensory areas. Importantly, representations in early visual areas (V1-V3AB)
163 showed a repulsive bias for all participants regardless of the specific decoding technique used and when
164 the decoder was instead trained on an independent orientation localizer (Figure S7-8). This suggests that
165 the repulsive bias is also found in the “sensory” code and is not specific to working memory
166 maintenance. This repulsive pattern held throughout the duration of the trial, suggesting it was not a
167 transient phenomenon (Figure S7A). In accordance with the large effects of the previous stimulus on
168 current trial representations, we observed above chance decoding for the identity of the previous
169 stimulus in 9/11 ROIs using the same TRs and decoding techniques as used for the current stimulus
170 (Figure S9).

171 We also examined how the precision of neural representations changed as a function of
172 stimulus history. In sharp contrast to behavior, σ_{circ} exhibited a monotonic trend such that neural
173 decoding was *least* precise when the previous stimulus was similar (Figure 2D, gray curve, see [Neural](#)
174 [Variance](#)). We quantified this difference in sensory uncertainty in a similar manner to the behavioral
175 data and found that variance in the sensory representations was significantly greater following a similar
176 stimulus ($<30^\circ$, $t(5)=13.33$, $p=.00002$, paired 1-tailed t-test, V3, Figure 2E). This pattern was significant
177 ($p<.01$) in 8/11 ROIs and the difference in precision was significant across ROIs ($t(10)=6.92$, $p=.00002$,
178 Figure S10A-B). The results did not change qualitatively when we utilized decoded uncertainty derived
179 directly from the posterior rather than the circular standard deviation of decoded responses³⁸ (Figure
180 S10C-D), or when we used other thresholds between 20° and 40° . The repulsion of sensory
181 representations and the corresponding reduction in decoding precision around the previous orientation
182 is consistent with neural adaptation where recently active units are attenuated, thus leading to lower
183 SNR responses in visual cortex.

184 [Encoder-Decoder Model](#)

185 We observed an attractive bias and low variability around the current stimulus feature in
186 behavior, and a repulsive bias and high variability around the current feature in the fMRI decoding data.

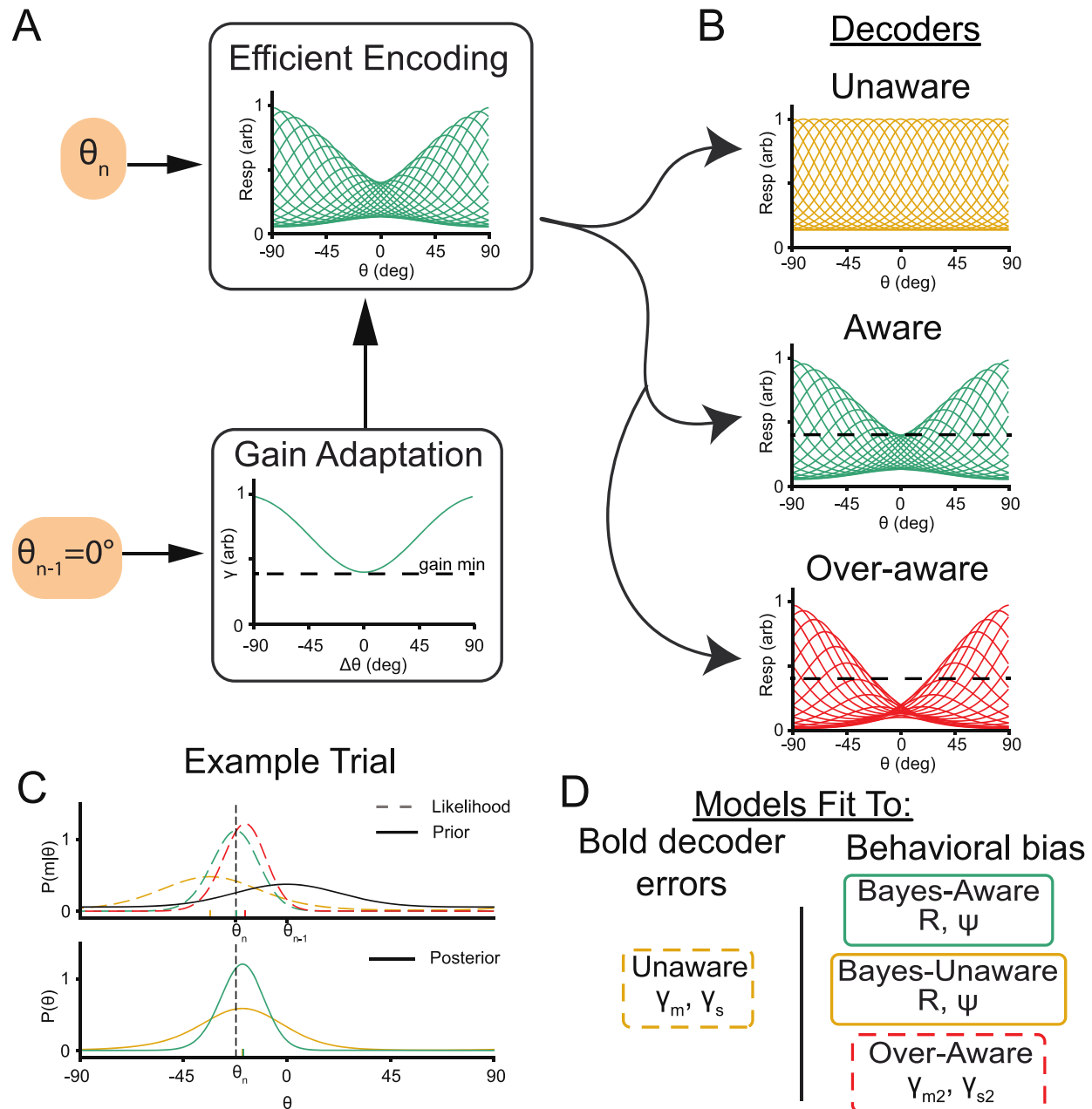


Figure 3 Encoder-Decoder model schematic. **A:** Encoding. Units with von Mises tuning curves encode incoming stimuli. The gain of individual units undergoes adaptation such that their activity is reduced as a function of their distance from the previous stimulus. **B:** Decoding. This activity is then read out using a scheme that assumes one of three adaptation profiles. The unaware decoder assumes no adaptation has taken place, the aware decoder assumes the true amount of adaptation while the over-aware decoder over-estimates the amount of adaptation (note center tuning curves dip lower than the minimum gain line from encoding). **C:** Example stimulus decoding. Top: The resulting likelihood function for the unaware readout (dotted yellow line) has its representation for the current trial ($\theta_n=-30^\circ$) biased away from the previous stimulus ($\theta_{n-1}=0^\circ$). The aware readout (dotted green line) is not biased, while the over-aware readout is biased towards the previous stimulus. These likelihood functions can be multiplied by a prior of stimulus contiguity (solid black line) to get a Bayesian posterior (bottom) where Bayes-unaware and Bayes-aware representations are shifted towards the previous stimulus. Tick marks indicate maximum likelihood or decoded orientation. **D:** Summary of models and free parameters being fit to both BOLD decoder errors and behavioral bias.

188 Thus, the patterns of bias and variability observed in the behavioral data are opposite to the patterns of
 189 bias and variability observed in visual cortex. To better understand these opposing effects, we reasoned
 190 that representations in early visual cortex do not directly drive behavior but instead are read out by later
 191 cortical regions that determine the correct response given the task³⁹⁻⁴¹. In this construction, the
 192 decoded orientations from visual cortex represent only the beginning of a complex information
 193 processing stream that, in our task, culminates with the participant making a speeded button press
 194 response. Thus, we devised a two-stage encoder-decoder model to describe observations in both early
 195 visual cortex and in behavior (see [modeling](#)).

196 The encoding stage of the model consists of a simulated population of orientation-selective
 197 neurons with von Mises tuning curves evenly tiling the feature space. The gain of these tuning curves
 198 undergoes adaptation such that units tuned to the previous stimulus (θ_{n-1}) will have their activity
 199 reduced on the current trial (Figure 3A). Note that these neurons are assumed to have Poisson firing
 200 rates and that their responses are noiseless while training the model.

201 The decoding stage reads out this activity using one of three strategies (Figure 3B). The *unaware*
 202 decoder assumes no adaptation has taken place and results in stimulus likelihoods $p(m|\theta)$ that are
 203 repelled from the previous stimulus (Figure 3C, yellow). This adaptation-naïve decoder is a previously
 204 hypothesized mechanism for behavioral adaptation³⁶ and likely what gives rise to the repulsive bias we
 205 observe in visual cortex using a fMRI decoder that is agnostic to stimulus history (Figure 2D).

	Fit To:	BOLD Decoder	Behavior		
Stage:		Unaware	Bayes unaware (Prior*unaware)	Bayes aware (Prior*aware)	Over-aware
Encoding	γ_m	X			
	γ_s	X			
Decoding	γ_{m2}	0	0	γ_m	X
	γ_{s2}	1	1	γ_s	X
Bayes	R	5	X	X	5
	ψ	N/A	X	X	N/A

Table 1 Cells correspond to parameters for proposed decoders, with 'X' indicating free parameters adjusted to fit empirical data. γ_m controls the amplitude and γ_s controls the width of gain adaptation (Figure 3A). These parameters are fit by minimizing the residual sum of squared errors between the unaware decoder and the BOLD decoder output. γ_{m2} and γ_{s2} are the assumed adaptation parameters at decoding. These terms are either set to assume no adaptation (unaware), match the true amount of adaptation (aware) or are free parameters adjusted to maximize the likelihood of responses (over-aware, Figure 3B). Last, R adjusts the average Poisson firing rate and ψ controls the variance of the prior distribution (Figure 3C). These parameters are adjusted for decoders using a Bayesian prior while R is set to the arbitrary value of 5 for non-Bayesian decoders (it has no effect on bias for non-Bayesian decoders). Increasing R increases the precision of the likelihood function and reduces the relative influence of the prior. Increasing ψ increases the range of $\Delta\theta$ over which the prior has an influence.

206 Alternatively, the *aware* decoder (Figure 3C, green) has perfect knowledge of the current state of
207 adaptation and can thus account for and ‘un-do’ biases introduced during encoding. Finally, the *over-*
208 *aware* decoder knows the identity of the previous stimulus but over-estimates the amount of gain
209 modulation that takes place, resulting in a net attraction to the previous stimulus (Figure 3C, red). We
210 additionally combined a formal prior based on temporal contiguity with the stimulus likelihood from the
211 previously described decoders³⁵. In our implementation, a Bayesian prior centered on the previous
212 stimulus (Figure 3C, black) is multiplied by the decoded likelihood to get a Bayesian posterior (Figure 3C,
213 bottom). We applied this prior of temporal contiguity to both the *aware* decoder as well as the *unaware*
214 decoder to test the importance of awareness at decoding. We did not apply a prior to the *over-aware*
215 model to balance the number of free parameters between the various decoders and to see if the *over-*
216 *aware* model could achieve attractive serial dependence without a Bayesian prior (Table 1). In total we
217 explored three separate decoder models: the *Bayes-unaware* and *Bayes-aware* models which apply a
218 prior to their respective likelihoods as well as the *over-aware* model which outputs the maximum
219 likelihood (Figure 3D).

220 For each participant, we fit the encoder-decoder model in two steps (Figure 3D). All model
221 fitting was performed using the same cross-validation groups as our BOLD decoder and each stage had
222 two free parameters that were fit using grid-search and gradient descent techniques. We first report
223 results from the encoding stage of the model. The gain applied at encoding was adjusted to minimize
224 the residual sum of squared errors (RSS) between the output of the *unaware* decoder and the residual
225 errors of our BOLD decoder. The *unaware* readout of the adapted encoding process (Figure 4A, yellow)
226 provided a good fit to the average decoding errors obtained with the BOLD decoder (Figure 4A, black
227 outline, $\rho=0.99$) and across individual participants (S11A, ranges: $\rho= [0.92,0.98]$). The *unaware* readout
228 provided a better fit to the outputs of our neural decoder than the presented orientation (Figure 4C,
229 $t(5)=5.94$, $p=.001$, paired one tailed t-test) because it captured a significant proportion of the variance in
230 decoding errors as a function of $\Delta\theta$ (Figure 4D, $t(5)=9.34$, $p=.0001$, one-tailed t-test).

231 We next considered three readout schemes of this adapted population to maximize the
232 likelihood of our behavioral responses (Figure 3D). The *Bayes-aware* decoder is consistent with previous
233 Bayesian accounts of serial dependence³⁸, but additionally asserts that Bayesian inference occurs after
234 encoding and that readout must account for adaptation. Alternatively, the *Bayes-unaware* decoder tests
235 whether this awareness is necessary to achieve attractive serial dependence. Both models were able to
236 achieve attractive biases that were positively correlated with average behavioral biases (Figure 4B), and
237 individual biases (Figure S11B-C), but the *Bayes-aware* model was significantly more likely given
238 participant responses (Figure 4E, $t(5)=12.8$, $p=5.3*10^{-5}$). We also considered the *over-aware* model to
239 determine if a mismatch between expected and true levels of adaptation can explain attractive serial
240 dependence without the need to invoke a formal Bayesian prior. This model also outperformed the
241 *Bayes-unaware* model (Figure 4E, $t(5)=3.69$, $p=.014$) but was not significantly different from the *Bayes-*
242 *aware* model ($p=0.18$, all t-tests paired, two-tailed). Finally, we examined the variance of the *unaware*
243 decoder as well as the three readout schemes fit to behavior (Table 1) to see if they were able to
244 reproduce patterns similar to the BOLD decoder and the behavioral responses, respectively. As model
245 coefficients were fit independent of observed variance, correspondence between model performance
246 and BOLD/behavioral data would provide convergent support for the best model. While the models
247 were trained using noiseless activity at encoding, we simulated responses using Poisson rates to induce
248 variability.

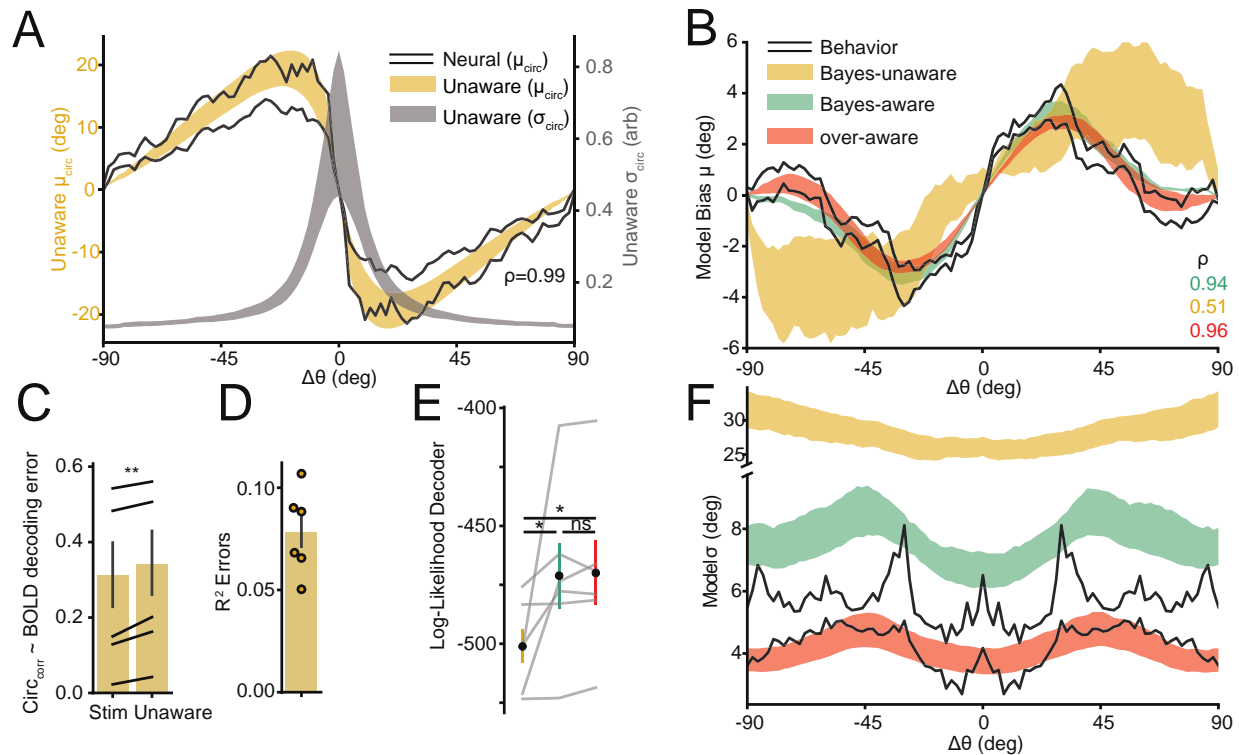


Figure 4 Model performance. **A,C,D**- Neural Decoder; **B,E,F**- Behavior. **A**: Unaware decoder (yellow) provides a good fit to neural bias (black outline). Decoded variance decreases monotonically with distance from previous stimulus. (\pm SEM across participants). **B**: Perceptual bias (black outline) was well fit by the Bayes-aware and over-aware models but not the Bayes-Unaware model (\pm SEM across participants). **C**: Unaware decoder provided a better fit to decoded responses than the presented orientation. **D**: The unaware decoder was able to explain a significant proportion of decoding error variance. **E**: Likelihood of observed responses for best fit model for each participant. Bayes-unaware significantly worse than other models. Same color scheme as **B**. **F**: Perceptual σ had a similar shape and magnitude to Bayes-aware and over-aware model fits. Bayes-unaware model output was much less precise and had a different form. Same legend as **B**. *, $p < .05$; **, $p < .01$; ***, $p < .001$

249 We simulated 1000 trials from each cross-validated fit and pooled the model outputs. We first
 250 confirmed that the variance of the *unaware* decoder was highest following small changes of $\Delta\theta$ (Figure
 251 4A, gray; Figure S12C $t(5)=3.4$, $p=.01$, paired 1-tailed t-test $<30^\circ$ vs $>30^\circ$) matching the output of our
 252 neural decoder (Figure 2G) and providing additional support for gain adaptation causing the observed
 253 repulsion in the fMRI data. Next, we examined the different behavioral decoders and found that,
 254 matching real behavioral responses, both the *Bayes-aware* and *over-aware* decoders were significantly
 255 more precise following small values of $\Delta\theta$ (Figure 4F; Figure S12C, Bayes-aware, $t(5)=-3.19$, $p=.012$, over-
 256 aware, $t(5)=-6.64$, $p=.0006$) while the *Bayes-unaware* decoder did not show this trend ($t(5)=-1.99$,
 257 $p=.052$). Notably, the overall magnitude of variance observed with the *Bayes-unaware* decoder was also
 258 much higher than that observed in the real behavioral data (Figure 4F) and thus provided a significantly
 259 worse fit relative to either of the aware decoders ($p < .005$, Figure S12A-B, paired t-tests comparing
 260 Jensen-Shannon divergence of error distributions). Together, the variance data provides additional
 261 evidence in favor of adaptation driving the repulsive biases that were observed in the BOLD data and
 262 awareness of the current state of adaptation being a requisite condition for attractive serial
 263 dependence.

264 Discussion

265 In this study, we sought to understand the neural underpinning of attractive serial dependence,
266 and how changes in tuning properties at encoding shape behavior. Based on previous behavioral and
267 neural studies, we expected to observe attractive biases in line with observed behavior and decoding
268 from early visual areas^{26,28,29}. Instead, we found that representations were significantly repelled from
269 the previous stimulus starting in primary visual cortex and continuing through IPS (Figure 2F). This
270 repulsion is consistent with bottom up adaptation beginning either at or before V1 and cascading up the
271 visual hierarchy^{11,42}. As repulsive biases are clearly in the opposite direction as behavioral biases, we
272 built a model to link these conflicting patterns. The critical new insight revealed by the model is that
273 only readout schemes that account for adaptation can explain attractive serial dependence. More
274 generally, our BOLD data strongly point against an early sensory or ‘perceptual’ account of serial
275 dependence and instead suggest that serial dependence is driven by post-perceptual or mnemonic
276 circuits^{34,43}.

277 Two previous studies have examined how sensory representations are shifted by serial
278 dependence. An fMRI study observed *attractive* biases in early visual areas that corresponded to
279 behavioral performance, but the study was limited as it only used two orthogonal orientations
280 ($\theta=\{45,135^\circ\}$)²⁸. As shown in the present results (Figures 1B, 2A) and in other studies^{26,29,32,44,45}, serial
281 dependence is absent at these offsets. In addition, the stimuli were rendered at low contrast and were
282 embedded in visual noise, making them difficult to accurately encode. Thus, the observed history bias
283 may be more akin to perceptual ‘priming’ as opposed to attractive serial dependence^{46,47}, as individuals
284 may have been able to detect stimuli faster and more reliably when they observed a similar stimulus in
285 the recent past (particularly on trials where they failed to accurately encode the near-threshold
286 orientation presented on the current trial). Further, as their decoder was trained with only two stimuli,
287 they could not build an explicit model of orientation representations. Thus, it is unclear if enhancements
288 at encoding correspond to a shift in orientation-selective information in voxel-activation patterns or
289 rather to a reduction in variability. More in line with the present experiment, a second study found that
290 population representations in FEF were *repelled* while saccades were attracted to the location of the
291 previous stimulus³². Perhaps because this effect was observed in a later visual area, the authors
292 explained their finding as a consequence of residual attentional shifts from the previous trial. Our
293 finding of repulsive biases as early as V1 is more consistent with bottom up adaptation as attention
294 effects tend to become more pronounced later in the visual hierarchy (and Papadimitriou and
295 colleagues also acknowledge this as an alternative mechanism)⁴⁸⁻⁵¹.

296 In line with classic accounts, adaptation in visual cortex should lead to a reduction in energy
297 usage during encoding¹⁰. However, our modeling results highlight the importance of an aware decoder,
298 which may offset adaptation-related efficiency gains. Instead the main advantage of adaptation may be
299 to decorrelate inputs, thus enhancing the discriminability of incoming stimuli^{9,10}. The resulting biases
300 may have little fitness cost relative to the advantage of being aware of stimulus changes potentially
301 signaling threat or food. Indeed an optimal processing stream may emphasize differences at encoding
302 and only favor stability once a stimulus has been selected by attention for more extensive post-
303 perceptual processing⁴³. This motif of pattern separation followed by pattern completion would not be
304 unique to adaptive visual processing. For example, similar mechanisms have been proposed as a critical
305 component of long term memory processing in the hippocampus and associative memory formation in

306 the fly mushroom body⁵². Thus, the biases introduced by adaptation may be beneficial in part because
307 they expand the dimensionality of the representational space.

308 In our model, we did not explicitly define how awareness of adaptation is implemented.
309 However, some representation of information about stimulus history appears to be a minimum
310 requirement. The identity of the previous stimulus for spatial position and angle has previously been
311 shown to be decodable from the spiking activity of single units in the frontal eye field (FEF) and large-
312 scale activity patterns in human EEG^{32,53}. We additionally demonstrate that information about the
313 previous trial is encoded in patterns of fMRI activity in human visual cortex (Figure S9). These signals
314 could potentially be represented concurrently with representations of the current stimulus in the same
315 populations of sensory neurons, but separating activity representing current and previous stimuli may
316 prove difficult under this scheme. An alternate, and potentially more appealing, account holds that
317 representations of stimulus history are maintained outside of early visual areas, consistent with findings
318 from mouse parietal cortex⁴. This anatomical segregation could disambiguate incoming sensory drive
319 from representations of stimulus history.

320 For the decoding stage of our model, we established that only readout schemes that are aware
321 of adaptation can explain attractive serial dependence. The *Bayes-aware* model is an extension of
322 previously proposed models that employ an explicit prior but that did not consider effects of adaptation
323 at encoding³⁵. In contrast, the *over-aware* model is a novel account that can achieve similar
324 performance without needing an explicit prior based on stimulus history. While model fit metrics did not
325 readily distinguish one of these two models as superior, the *over-aware* model may prove to be more
326 flexible. For instance, one of our fMRI participants showed significant repulsion from far stimuli, an
327 observation also reported by others^{29,31}. While the *over-aware* model can fit this repulsive regime, the
328 *Bayes-aware* model is incapable of generating repulsive patterns (compare models fits for subj #3,
329 Figure S11). This limitation of a purely Bayesian account of serial dependence is also observable in prior
330 work (Figure 6B in³⁵).

331 The *over-aware* – or more generally a “flexibly-aware” – decoder may also account for
332 phenomena not covered in the present study. While behavioral (repulsive) adaptation is assumed to
333 result from an unaware decoder^{8,36}, the magnitude of neural adaptation may be much larger than the
334 resulting behavioral repulsion observed^{12,19,23}. Thus, behavioral adaptation may arise when adaptation
335 outweighs awareness (an ‘*under-aware*’ decoder) which could arise in paradigms where inducing stimuli
336 are task irrelevant and presented for long periods of time^{19,54}. By contrast, *over-aware* decoders may
337 arise in laboratory paradigms that involve attending to and holding in memory a weak stimulus^{26,27}.

338 In this study, we extended previous descriptions of serial dependence by quantifying how both
339 bias and variance are shaped by stimulus history. We report a robust pattern of perception being most
340 precise following small changes in successive stimulus features (Figure 1C-D, 2A-B). This relationship
341 violates a perceptual ‘law’ proposing that bias is inversely proportional to the derivative of
342 discrimination thresholds⁵⁵. This ‘law’ would assert that our attractive bias should come with a less
343 precise representation following small changes (or a repulsive bias to account for our enhanced
344 precision). We argue that serial dependence is not ‘violating’ this law, but rather believe this is further
345 evidence for serial dependence being a post-perceptual phenomenon. Neural representations exhibit
346 repulsive biases, expanding the perceptual space and allowing greater discriminability. When these

347 representations are read out by an aware decoder, the bias is undone but the enhanced discriminability
348 remains (Figure 4F).

349 Acknowledgements

350 Funded by NEI R01-EY025872 to JTS, and NIMH Training Grant in Cognitive Neuroscience (T32-
351 MH020002) to TS. Thanks to Chaipat Chunharas for critical discussions in experimental design and
352 assistance with scanning and to Anika Jollorina and Shuangquan Feng for assistance with behavioral data
353 collection. Code for Bayesian decoder adapted from code provided by Ruben van Bergen. Thanks to
354 Marcelo Mattar for helpful comments on our model and to Margaret Henderson, Sunyoung Park, and
355 Kirsten Adam for thoughtful comments on earlier versions of the manuscript.

Works Cited

1. Dong, D. W. & Atick, J. J. Statistics of natural time-varying images. *Netw. Comput. Neural Syst.* **6**, 345–358 (1995).
2. Felsen, G., Touryan, J. & Dan, Y. Contextual modulation of orientation tuning contributes to efficient processing of natural stimuli. *Netw. Comput. Neural Syst.* **16**, 139–149 (2005).
3. Girshick, A. R., Landy, M. S. & Simoncelli, E. P. Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nat. Neurosci.* **14**, 926–932 (2011).
4. Akrami, A., Kopec, C. D., Diamond, M. E. & Brody, C. D. Posterior parietal cortex represents sensory history and mediates its effects on behaviour. *Nature* **554**, 368–372 (2018).
5. Ashourian, P. & Loewenstein, Y. Bayesian Inference Underlies the Contraction Bias in Delayed Comparison Tasks. *PLoS ONE* **6**, e19551 (2011).
6. Wei, X.-X. & Stocker, A. A. A Bayesian observer model constrained by efficient coding can explain ‘anti-Bayesian’ percepts. *Nat. Neurosci.* **18**, 1509–1517 (2015).
7. Barlow, H. B. Possible Principles Underlying the Transformations of Sensory Messages. in *Sensory Communication* (ed. Rosenblith, W. A.) 216–234 (The MIT Press, 2012).
doi:10.7551/mitpress/9780262518420.003.0013.
8. Benucci, A., Saleem, A. B. & Carandini, M. Adaptation maintains population homeostasis in primary visual cortex. *Nat. Neurosci.* **16**, 724–729 (2013).
9. Clifford, C. W. G. *et al.* Visual adaptation: Neural, psychological and computational aspects. *Vision Res.* **47**, 3125–3131 (2007).
10. Clifford, C. W. G., Wenderoth, P. & Spehar, B. A functional angle on some after-effects in cortical vision. *Proc. R. Soc. Lond. B Biol. Sci.* **267**, 1705–1710 (2000).
11. Dhruv, N. T. & Carandini, M. Cascaded Effects of Spatial Adaptation in the Early Visual System. *Neuron* **81**, 529–535 (2014).

12. Dragoi, V., Rivadulla, C. & Sur, M. Foci of orientation plasticity in visual cortex. *Nature* **411**, 80–86 (2001).
13. Dragoi, V., Sharma, J. & Sur, M. Adaptation-Induced Plasticity of Orientation Tuning in Adult Visual Cortex. *Neuron* **28**, 287–298 (2000).
14. Durant, S., Clifford, C. W. G., Crowder, N. A., Price, N. S. C. & Ibbotson, M. R. Characterizing contrast adaptation in a population of cat primary visual cortical neurons using Fisher information. *J. Opt. Soc. Am. A* **24**, 1529 (2007).
15. Kohn, A. & Movshon, J. A. Adaptation changes the direction tuning of macaque MT neurons. *Nat. Neurosci.* **7**, 764–772 (2004).
16. Gardner, J. L. *et al.* Contrast Adaptation and Representation in Human Early Visual Cortex. 14 (2005).
17. Patterson, C. A., Wissig, S. C. & Kohn, A. Adaptation Disrupts Motion Integration in the Primate Dorsal Stream. *Neuron* **81**, 674–686 (2014).
18. Dekel, R. & Sagi, D. Tilt aftereffect due to adaptation to natural stimuli. *Vision Res.* **117**, 91–99 (2015).
19. He, S. & MacLeod, D. I. A. Orientation-selective adaptation and tilt after-effect from invisible patterns. *Nature* **411**, 473–476 (2001).
20. Moradi, F., Koch, C. & Shimojo, S. Face Adaptation Depends on Seeing the Face. *Neuron* **45**, 169–175 (2005).
21. Abbonizio, G., Clifford, C. & Langley, K. Contrast adaptation may enhance contrast discrimination. *Spat. Vis.* **16**, 45–58 (2002).
22. Clifford, C. W. G., Wyatt, A. M., Arnold, D. H., Smith, S. T. & Wenderoth, P. Orthogonal adaptation improves orientation discrimination. *Vision Res.* **41**, 151–159 (2001).

23. Jin, D. Z., Dragoi, V., Sur, M. & Seung, H. S. Tilt Aftereffect and Adaptation-Induced Changes in Orientation Tuning in Visual Cortex. *J. Neurophysiol.* **94**, 4038–4050 (2005).
24. Phinney, R. E., Bowd, C. & Patterson, R. Direction-selective Coding of Stereoscopic (Cyclopean) Motion. *Vision Res.* **37**, 865–869 (1997).
25. Cicchini, G. M., Mikellidou, K. & Burr, D. Serial dependencies act directly on perception. *J. Vis.* **17**, 6 (2017).
26. Fischer, J. & Whitney, D. Serial dependence in visual perception. *Nat. Neurosci.* **17**, 738–743 (2014).
27. Manassi, M., Liberman, A., Kosovicheva, A., Zhang, K. & Whitney, D. Serial dependence in position occurs at the time of perception. *Psychon. Bull. Rev.* **25**, 2245–2253 (2018).
28. St. John-Saaltink, E., Kok, P., Lau, H. C. & de Lange, F. P. Serial Dependence in Perceptual Decisions Is Reflected in Activity Patterns in Primary Visual Cortex. *J. Neurosci.* **36**, 6186–6192 (2016).
29. Fritsche, M., Mostert, P. & de Lange, F. P. Opposite Effects of Recent History on Perception and Decision. *Curr. Biol.* **27**, 590–595 (2017).
30. Cicchini, G. M., Mikellidou, K. & Burr, D. C. The functional role of serial dependence. *Proc. R. Soc. B Biol. Sci.* **285**, 20181722 (2018).
31. Bliss, D. P., Sun, J. J. & D’Esposito, M. Serial dependence is absent at the time of perception but increases in visual working memory. *Sci. Rep.* **7**, 14739 (2017).
32. Papadimitriou, C., White, R. L. & Snyder, L. H. Ghosts in the Machine II: Neural Correlates of Memory Interference from the Previous Trial. *Cereb. Cortex* bhw106 (2016)
doi:10.1093/cercor/bhw106.
33. Papadimitriou, C., Ferdoash, A. & Snyder, L. H. Ghosts in the machine: memory interference from the previous trial. *J. Neurophysiol.* **113**, 567–577 (2015).
34. Barbosa, J. *et al.* Interplay between persistent activity and activity-silent dynamics in the prefrontal cortex underlies serial biases in working memory. *Nat. Neurosci.* **23**, 1016–1024 (2020).

35. van Bergen, R. S. & Jehee, J. F. M. Probabilistic Representation in Human Visual Cortex Reflects Uncertainty in Serial Decisions. *J. Neurosci.* **39**, 8164–8176 (2019).
36. Seriès, P., Stocker, A. A. & Simoncelli, E. P. Is the Homunculus “Aware” of Sensory Adaptation? *Neural Comput.* **21**, 3271–3304 (2009).
37. Cicchini, G. M. & Burr, D. C. Serial effects are optimal. *Behav. Brain Sci.* **41**, e229 (2018).
38. van Bergen, R. S., Ji Ma, W., Pratte, M. S. & Jehee, J. F. M. Sensory uncertainty decoded from visual cortex predicts behavior. *Nat. Neurosci.* **18**, 1728–1730 (2015).
39. Crick, F. & Koch, C. Are we aware of neural activity in primary visual cortex? *Nature* **375**, 121–123 (1995).
40. Gold, J. I. & Shadlen, M. N. The Neural Basis of Decision Making. *Annu. Rev. Neurosci.* **30**, 535–574 (2007).
41. Grunewald, A., Bradley, D. C. & Andersen, R. A. Neural Correlates of Structure-from-Motion Perception in Macaque V1 and MT. *J. Neurosci.* **22**, 6195–6207 (2002).
42. Patterson, C. A., Duijnhouwer, J., Wissig, S. C., Krekelberg, B. & Kohn, A. Similar adaptation effects in primary visual cortex and area MT of the macaque monkey under matched stimulus conditions. *J. Neurophysiol.* **111**, 1203–1213 (2014).
43. Pascucci, D. *et al.* Laws of concatenated perception: Vision goes for novelty, decisions for perseverance. *PLOS Biol.* **17**, e3000144 (2019).
44. Kiyonaga, A., Scimeca, J. M., Bliss, D. P. & Whitney, D. Serial Dependence across Perception, Attention, and Memory. *Trends Cogn. Sci.* **21**, 493–497 (2017).
45. Fischer, C. *et al.* Context information supports serial dependence of multiple visual objects across memory episodes. *Nat. Commun.* **11**, 1932 (2020).
46. Bar, M. & Biederman, I. Subliminal Visual Priming. *Psychol. Sci.* **9**, 6 (1998).
47. Tulving, E. & Schacter, D. Priming and human memory systems. *Science* **247**, 301–306 (1990).

48. Kastner, S. Mechanisms of Directed Attention in the Human Extrastriate Cortex as Revealed by Functional MRI. *Science* **282**, 108–111 (1998).
49. Luck, S. J., Chelazzi, L., Hillyard, S. A. & Desimone, R. Neural Mechanisms of Spatial Selective Attention in Areas V1, V2, and V4 of Macaque Visual Cortex. *J. Neurophysiol.* **77**, 24–42 (1997).
50. Sprague, T. C. & Serences, J. T. Attention modulates spatial priority maps in the human occipital, parietal and frontal cortices. *Nat. Neurosci.* **16**, 1879–1887 (2013).
51. Treue, S. Visual attention: the where, what, how and why of saliency. *Curr. Opin. Neurobiol.* **13**, 428–432 (2003).
52. Cayco-Gajic, N. A. & Silver, R. A. Re-evaluating Circuit Mechanisms Underlying Pattern Separation. *Neuron* **101**, 584–602 (2019).
53. Bae, G.-Y. & Luck, S. J. Reactivation of Previous Experiences in a Working Memory Task. *Psychol. Sci.* **30**, 587–595 (2019).
54. Gibson, J. J. & Radner, M. Adaptation, after-effect and contrast in the perception of tilted lines. I. Quantitative studies. *J. Exp. Psychol.* **20**, 453–467 (1937).
55. Wei, X.-X. & Stocker, A. A. Lawful relation between perceptual bias and discriminability. *Proc. Natl. Acad. Sci.* **114**, 10244–10249 (2017).
56. Roth, Z. N., Heeger, D. J. & Merriam, E. P. Stimulus vignetting and orientation selectivity in human visual cortex. *eLife* **7**, e37241 (2018).
57. SciPy 1.0 Contributors *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).
58. Stein, H. *et al.* Reduced serial dependence suggests deficits in synaptic potentiation in anti-NMDAR encephalitis and schizophrenia. *Nat. Commun.* **11**, 4250 (2020).
59. Averbeck, B. B., Latham, P. E. & Pouget, A. Neural correlations, population coding and computation. *Nat. Rev. Neurosci.* **7**, 358–366 (2006).

356 Methods

357 Participants

358 Behavioral study: 56 participants (male and female) were drawn from a subject pool of primarily
359 undergraduate students at UC San Diego. All subjects gave written consent to participate in the study in
360 accordance with the UC San Diego IRB, and were compensated either monetarily or with class credit. Of
361 these 56 participants, 9 were removed from further analysis for completing less than 200 trials (2) or
362 getting less than 60% of trials correct (7). We included the remaining 47 participants who completed on
363 average 421 trials, range: [204, 988], in our lab over the course of 1 to 3 sessions.

364 fMRI study: 6 participants (3 female, mean age 24.6 ± 0.92) participated in four, 2-hour scanning
365 sessions. Each subject completed between 748 and 884 trials (mean 838.7). For two participants, one
366 session had to be repeated due to technical difficulties that arose during scanning.

367 Behavioral Discrimination Task

368 Participants in the behavior-only study completed the task on a desktop computer in a sound
369 attenuated room. Subjects were seated with a chin rest to stabilize viewing 50 cm from a 39 by 29 cm
370 CRT monitor (1600x1200 px) with a visual angle of 42.6° (screen width). Each trial consisted of a full-field
371 oriented grating (1000 ms) which had to be remembered across a delay period (3,500 ms) before a test.
372 At test, the participant judged whether a line was slightly clockwise (CW) or counter-clockwise (CCW)
373 relative to the remembered orientation (max response time window: 3,000ms, Figure 1A). The oriented
374 grating consisted of a sine wave grating (spatial frequency $1.73 \text{ cycles}/^\circ$, 0.8 Michelson contrast)
375 multiplied by a 'donut' mask (outer diameter $\phi=24.3^\circ$, inner $\phi=1.73^\circ$). The stimulus was then convolved
376 with a 2D Gaussian filter (1.16° kernel, $SD = 0.58^\circ$) to minimize edge artifacts⁵⁶. Phase and orientation
377 were randomized across trials, and the stimulus was phase-reversed every 250ms. After the offset of the
378 oriented grating, a mask of filtered noise was presented for 500ms. The mask was generated by band
379 passing white noise [low 0.22 , high $0.87 \text{ cycles}/^\circ$], multiplying by the same donut mask, and convolving
380 with a 2D Gaussian filter (0.27° kernel, $SD = 0.11^\circ$). The mask was phase reversed once after 250 ms. A
381 black fixation point (diameter $.578^\circ$) was displayed throughout the extent of the block and turned white
382 for 500 ms prior to stimulus onset on each trial. The probe was a white line (width 0.03° , length 24.3°)
383 masked by the same donut. Subjects indicated whether the probe line was CW or CCW from the
384 remembered orientation by pressing one of two buttons ('Q', 'P') with their left and right pointer
385 fingers. The next trial started after a 1000ms inter trial interval (ITI). For some behavioral participants
386 ($n=9$) delay and ITI were varied between 0.5-7.5s without notable effects on performance.

387 First, subjects completed a training block to ensure that they understood the task. Next, they
388 completed a block of trials where difficulty was adjusted by changing the probe offset ($\delta\theta$) between the
389 stimulus and probe to achieve 70% accuracy. This $\delta\theta$ was used in subsequent blocks and was adjusted
390 on a per-block basis to keep performance at approximately 70%. Participants completed an average of
391 5.76 ± 0.24 blocks [min = 3, max = 9]. Some participants completed the task with slight variations in the
392 distribution and sequence of orientations presented. For completeness we include those details here.
393 Note, however, we additionally report a set of control analyses in which we repeat all of our main
394 analyses excluding blocks with binned stimuli and find no relevant difference in behavior. For most
395 participants, stimuli were pseudo-randomly distributed across the entire 180° space such that they were
396 uniformly distributed across blocks of 64 trials ($n=25$). However, some participants saw stimuli that were
397 binned (with some jitter) every 22.5° to purposefully avoid cardinal and oblique orientations (11.25° ,

398 33.75°, 56.25°, etc.) and the trial sequence was ordered so that a near oblique orientation was always
399 followed by a near cardinal orientation (n=7). This was implemented to maximize our ability to observe
400 serial dependencies in our binary response data as it is typically strongest around orientation changes of
401 20° and is more pronounced around oblique orientations³⁷. The remaining participants completed both
402 blocks with uniform and blocks with binned stimuli (n=14). All participants were interviewed after the
403 study and reported that stimuli were non-predictable and that all orientations felt equally likely. For our
404 main analysis we include all trials from all participants, irrespective of whether they participated in
405 uniform blocks, binned blocks, or both.

406 [fMRI Discrimination Task](#)

407 In the scanner, participants completed the behavioral task outlined above with slight
408 modifications. fMRI participants completed the task using a fiber-optic button box while viewing stimuli
409 through a mirror projected onto a screen mounted inside of the bore. The screen was 24 by 18 cm and
410 was viewed at a distance of 47 cm (width: 28.6° visual angle; 1024x768 px native resolution). The
411 stimulus timing was the same except that the sample-to-probe delay period was either 5, 7 or 9 s and
412 the ITIs were uniformly spaced between 5s and 9s and shuffled pseudo-randomly on each run of 17
413 trials. The oriented gratings had a spatial frequency of 1.27 cycles/°, outer $\varnothing=21.2^\circ$, inner $\varnothing=2.37^\circ$ and
414 were smoothed by a Gaussian filter (0.79° kernel, sd=0.79°). The noise patch (SF low 0.16, high 0.63
415 cycles/°) was also smoothed by a Gaussian filter (0.29° kernel, sd=0.11°). The probe stimulus was a white
416 line (width = 0.03°).

417 fMRI participants completed 44-52 blocks of 17 trials spread across 4, two-hour scanning
418 sessions for a total of 748-884 trials. As in the behavior-only task described above, 4 out of 6 fMRI
419 subjects had some blocks of trials where the stimuli were binned in 22.5° increments and ordered in a
420 non-independent manner (21-24 blocks/participant). However, all of the fMRI subjects also participated
421 in blocks with a uniform distribution of orientations across the entire 180° space (24-52
422 blocks/participant). For our main analysis we include all trials from all participants. However, as with the
423 behavioral analyses, we also report control analyses in which we repeat all of our main analyses
424 excluding blocks with non-random stimuli.

425 [fMRI Localizer Task](#)

426 Interleaved between the main task blocks, participants completed an independent localizer task
427 used for [voxel selection](#) where they were presented with a sequence of grating stimuli at different
428 orientations. Stimuli had a pseudo-randomly determined orientation that either matched the spatial
429 location occupied by the *donut* stimuli used in our main task (outer diameter $\varnothing=21.2^\circ$, inner diameter
430 $\varnothing=2.37^\circ$) or were a smaller foveal oriented Gabor corresponding to the 'hole' in the *donut* stimuli
431 (diameter $\varnothing=2.37^\circ$). Participants were instructed to attend to one of three features orthogonal to
432 orientation depending on the block: detect a contrast change across the entire stimulus, detect a small
433 grey blob appearing over part of the stimulus, or detect a small change in contrast at the fixation point.
434 Each stimulus was presented for 6000 ms and was separated by an ITI ranging from 3-8s.

435 [Response Bias](#)

436 Each trial consisted of a stimulus and a probe separated by a probe offset ($\delta\theta$) that was either
437 positive (probe is CW of stimulus) or negative. Participants judged whether the probe was CW or CCW
438 relative to the remembered orientation by making a binary response. To quantify the precision and the
439 response bias, we fit participant responses with a Gaussian cumulative density function with parameters

440 μ and σ corresponding to the *bias (mean)* and *standard deviation* of the distribution. The likelihood of a
441 given distribution was determined by the area under the curve (AUC) of the distribution of CW (CCW)
442 offsets between the stimulus and the probe ($\delta\theta$) on trials where the participant responded CW (CCW;
443 see Figure S1). In extreme cases, a very low standard deviation (σ) value with no bias would mean that
444 all $\delta\theta$ would lie outside the distribution and the participant would get every trial correct (Figure S1A). A
445 high negative bias (μ) value would mean that $\delta\theta$ would always lie CW relative to the distribution and the
446 participant would respond CW on every trial (Figure S1B). The best fitting parameters were found using
447 a bounded minimization algorithm (limited memory BFGS) on the negative log likelihood of the resulting
448 responses (excluded the small number of trials without a response) given the generated distribution⁵⁷.
449 We included a constant 25% guess rate in all model fits to ensure the likelihood of any response could
450 never be 0 (critical for later modelling). While this was critical to fitting our model to raw data, the
451 specific choice had no qualitative effect on our behavioral findings besides making the σ values smaller
452 compared to having a 0% guess rate. By having a constant guess rate rather than varying it as a free
453 parameter we were able to directly compare σ values across participants as a measure of performance.
454 Realistic model parameters and the effects of bias on response likelihood are also demonstrated (Figure
455 S1 C-D).

456 Serial Dependence

457 To quantify the dependence of responses on previous stimuli, we analyzed response bias and
458 variance as a function of the difference in orientation between the previous and current orientation
459 ($\Delta\theta = \theta_{n-1} - \theta_n$). We performed this analysis using a sliding window of 16° . To improve power, we
460 ‘folded’ our response data such that, when examining bias at 30° we included values from $22^\circ - 38^\circ$ as
461 well as responses from $-22^\circ - (-38^\circ)$ by inverting both the responses and probe offsets ($\delta\theta$) for the
462 negative values of $\Delta\theta$. This procedure removes any systematic responses biases (e.g., favoring CW
463 responses) and, as a result, the figures presenting serial dependence have rotational symmetry across
464 the origin^{34,58}.

465 We additionally fit a Derivative of Gaussian (DoG) function to parameterize the bias of
466 participant responses. The DoG function is parameterized with an amplitude A and width w

$$467 \quad y = xAwce^{-(wx)^2} \quad [1]$$

468 where $c = \sqrt{2e}$ is a normalization constant. For the purpose of fitting to our participant responses, x is
469 $\Delta\theta$ and y corresponds to μ in our response model. For each participant we adjusted three parameters: A ,
470 w , and σ to maximize the likelihood of participant responses. We report the magnitude of our fits as well
471 as the resulting full width at half max (FWHM) estimated numerically.

472 Response Precision

473 In addition to quantifying how responses were biased as a function of stimulus history, we also
474 estimated how precise responses were depending on their unsigned distance from the previous stimulus
475 ($|\Delta\theta|$). We used the same ‘folding’ procedure described in the previous section and only included trials
476 on the right half of our bias/variance plots (eg. Figure 1C, $\Delta\theta > 0$) to avoid double counting trials. Values
477 from the bin with more samples (typically ‘far’) were resampled (31 repetitions) without replacement
478 with the number of samples in the smaller bin and the median chosen to control for sample number
479 differences.

480 Scanning

481 fMRI task images were acquired over the course of four 2-hour sessions for each participant in a
482 General Electric Discovery MR750 3.0T scanner at the UC San Diego Keck Center for Functional Magnetic
483 Resonance Imaging. Functional echo-planar imaging (EPI) data were acquired using a Nova Medical 32-
484 channel head coil (NMSC075-32- 3GE-MR750) and the Stanford Simultaneous Multi-Slice (SMS) EPI
485 sequence (MUX EPI), with a multiband factor of 8 and 9 axial slices per band (total slices 72; 2-mm³
486 isotropic; 0-mm gap; matrix 104 x 104; field of view 20.8 cm; TR/TE 800/35 ms; flip angle 52°; in-plane
487 acceleration 1). Image reconstruction and un-aliasing was performed on cloud-based servers using
488 reconstruction code from the Center for Neural Imaging at Stanford. The initial 16 repetition times (TRs)
489 collected at sequence onset served as reference images required for the transformation from k-space to
490 the image space. Two 17s runs traversing k-space using forward and reverse phase-encoding directions
491 were collected in the middle of each scanning session and were used to correct for distortions in EPI
492 sequences using FSL top-up (FMRIB Software Library) for all runs in that session (Andersson et al. 2013,
493 Jenkinson et al. 2012). Reconstructed data was motion corrected and aligned to a common image. Voxel
494 data from each run was de-trended (8TR filter) and z-scored.

495 We also acquired one additional high-resolution anatomical scan for each subject (1 x 1 x 1-mm³
496 voxel size; TR 8,136 ms; TE 3,172 ms; flip angle 8°; 172 slices; 1-mm slice gap; 256x192-cm matrix size)
497 during a separate retinotopic mapping session using an Invivo eight-channel head coil. This scan
498 produced higher quality contrast between gray and white matter and was used for segmentation,
499 flattening, and visualizing retinotopic mapping data. The functional retinotopic mapping scanning was
500 collected using the 32-channel coil described above and featured runs where participants viewed
501 checkerboard gratings while responding to an orthogonal feature (transient contrast changes). Separate
502 runs featured alternating vertical and horizontal bowtie stimuli; rotating wedges; and an expanding
503 donut to generate retinotopic maps of the visual meridian, polar angle, and eccentricity respectively
504 (see Sprague and Serences, 2013). These images were processed using FreeSurfer and FSL functions and
505 visual regions of interest (ROI) were manually drawn on surface reconstructions (for areas: V1-V3, V3AB,
506 hV4, IPS0-IPS2, VO, LO, and TO).

507 Voxel Selection

508 To include only voxels that showed selectivity for the location of the oriented grating stimulus
509 used in our main experimental task, we used responses evoked during the independent *localizer* task
510 (see [fMRI Localizer Task](#)). For all analysis we used TRs 5-11 (4-8.8s) following stimulus onset. First, voxels
511 were selected based on their response to the spatial location of the grating stimulus by performing a t-
512 test on the responses of each voxel evoked by the donut and the donut-hole stimuli, selecting the 50%
513 of the voxels most selective to the donut for a given ROI. Of the voxels that passed this cutoff, we then
514 performed an ANOVA across 10° orientation bins and selected the 50% of voxels with the largest F-score
515 thus retaining ~25% of the initial voxel pool. These selected voxels were used in all main analysis.

516 Orientation Decoding (BNC)

517 We performed orientation decoding on BOLD activation patterns using a sliding temporal
518 window of 4 TRs. For most analysis we focused on a 3.2s (4 TR) window centered 6.4 s after stimulus
519 presentation. The Bayesian Noise Correlation (BNC) decoder assumes voxels are composed of
520 populations of neurons with tuning functions centered on one of 8 orientations evenly tiling the 180°
521 space. The response of population *i* to stimulus θ is given by:

$$c_i(\theta) = \max(0, \cos^5(\theta - \varphi_i)) \quad [2]$$

522 where φ_i is the center of the tuning function. The response of voxel j is defined as a weighted sum of
523 these hypothetical populations:

$$B_j = \sum_i^8 c_i w_i \quad [3]$$

524 Or in matrix notation,

$$B = CW \quad [4]$$

525

526 Where B (*trial x voxel*) is the resulting BOLD activity, C (*trial x channel*) is the hypothetical population
527 response, and W (*channel x voxel*) is the weight matrix. The weight matrix W is estimated as:

$$\hat{W} = C^{-1}B \quad [5]$$

528 where C^{-1} (*channel x trial*) is the pseudo-inverse of C (implemented using the NumPy `pinv` function).
529 This model was used to generate a linear estimate of voxel responses. The resulting residuals
530 correspond to voxel noise.

$$\hat{B} = C\hat{W} \quad [6]$$

531

$$B_{Noise} = B - \hat{B} \quad [7]$$

532 Noise correlations are known to contribute to observed activity and can be detrimental to our
533 resulting decoding capabilities⁵⁹. To reduce the impact of noise correlations across similarly tuned
534 populations, we implemented a Bayesian decoder that explicitly models these correlations³⁸. Briefly, we
535 modeled noise as coming from 3 distinct components: *global noise* shared across all voxels, *channel*
536 *noise* shared across neurons with similar tuning, and *voxel noise* explaining residual fluctuations in
537 individual voxels (see³⁸ for more details). The magnitude of these noise sources was estimated through
538 maximizing the likelihood of the observed residuals using a multivariate Gaussian defined by (number of
539 voxels) + (1 global) + (1 channel) parameters. After estimating noise sources, we could estimate the
540 posterior probability distribution given our fit weights \hat{W} and noise parameters $\hat{\Omega}$:

$$P(\theta|B; \hat{W}, \hat{\Omega}) = \frac{P(\theta|B; \hat{W}, \hat{\Omega})p(\theta)}{\int P(\theta|B; \hat{W}, \hat{\Omega})p(\theta) d\theta} \quad [8]$$

541 For each trial we then selected the θ most likely to have given rise to response B given
542 \hat{W} and $\hat{\Omega}$ as our decoded orientation and used its vector length as a proxy for model certainty. The
543 encoding and noise parameters of our model were fit to a subset of data and used to estimate
544 responses for held-out trials of the task data. We used leave-1-block-out cross-validation where each
545 block was a set of 4 consecutive runs (64 trials). These blocks featured orientations that were evenly
546 distributed across the entire 180° space to ensure a balanced training set. We performed additional
547 analysis training a model on the localizer task and testing on the memory task as well as cross-validating

548 within the localizer task. These models had lower SNR than models trained on the task but showed
549 qualitatively similar results as our task trained neural decoder.

550

551 Orientation Decoding (IEM)

552 For some analysis we additionally include the outputs of an Inverted Encoding Model (IEM). The
553 IEM uses the same encoding model as the BNC decoder (eq. [2-5]) but does not generate a specific
554 model of noise covariance. We instead inverted our estimated weight matrix (\widehat{W}) to estimate the
555 channel response on held out trials, $\widehat{C} = B\widehat{W}^{-1}$, where \widehat{W}^{-1} is the pseudo-inverse of \widehat{W} . The circular
556 mean of \widehat{C} was taken as the orientation estimate.

557 Neural Bias

558 To quantify how BOLD representations were biased by sensory history we computed the circular
559 mean of decoding errors ($\theta_{\text{error}} = \text{wrap}(\theta_{\text{decode}} - \theta_{\text{stim}})$):

$$\mu_{\text{circ}} = \text{angle}(\vec{R}), \quad [9]$$

560

$$\vec{R} = \frac{1}{n\text{Trials}} \sum_{k=0}^{n\text{Trials}} e^{i\theta_{\text{error}}^k}. \quad [10]$$

561

562 We estimated this bias using the same 16° sliding window as a function of $\Delta\theta$ used for visualizing
563 response bias from participant responses. We additionally quantified the magnitude of the bias in
564 decoding errors by fitting a DoG function to the raw decoding errors by minimizing the residual sum of
565 squares (RSS) and reporting the amplitude term.

566 Neural Variance

567 To quantify the variance of decoded orientations from visual areas, we computed the circular
568 standard deviation:

$$\sigma_{\text{circ}} = \sqrt{-2 \ln |\vec{R}|}. \quad [11]$$

569

570 This was estimated on both binned decoding errors (eq. 10), or on the single trial posterior estimate
571 from our orientation decoder (eq. 8). This was visualized using the same sliding window analysis as well
572 as in reference to whether it was close or far from the previous stimulus. Both pooled and single trial
573 estimates are reported and give similar results.

574 Modeling

575 We sought to develop a model that could explain both neural and behavioral biases as a
576 function of stimulus history. For the fMRI data, we focused on explaining changes in encoding that could
577 lead to the observed biases in the output of the BOLD decoder that was specifically designed to be
578 ‘unaware’ of stimulus history. To explain the behavioral data, we assumed that a decoder would receive
579 inputs from the same population of sensory neurons that we measured with fMRI and that the decoder
580 would read out this information in a manner that gives rise to attractive serial dependence. We
581 considered readout models that were either *unaware*, *aware*, or *over-aware* of adaptation and

582 additionally applied a Bayesian inference stage, which integrates prior expectations of temporal
 583 stability, to the *unaware* and *aware* decoders³⁵. We then compared performance between these
 584 competing models to see which could best explain our behavioral data.

585 Our full models consisted of two stages: an encoding stage where the gain of artificial neurons
 586 was changed as a function of the previous stimulus (adaptation) and a decoding stage where the
 587 readout from this adapted population was modified. The encoding population consisted of 100 neurons
 588 with von Mises tuning curves evenly tiling the 180° space. The expected unadapted population response
 589 is:

$$Resp_N(\theta_n) = R \gamma_N e^{\kappa \cos(\Phi - \theta_n) - 1} \quad [12]$$

590 Where γ_N is the scalar 1 for constant gain without adaptation, Φ is the vector of tuning curve
 591 centers, θ_n is the orientation of the current stimulus, $\kappa=1.0$ is a constant controlling tuning width, and R
 592 is a general gain factor driving the average firing rate. We implemented sensory adaptation by adjusting
 593 the gain of tuning curves relative to the identity of the previous stimulus, θ_{n-1} (Figure 3A, *Gain*
 594 *Adaptation*):

$$\gamma_A(\theta_{n-1}) = \gamma_N - rect(\gamma_m \cos^3(\gamma_s(\Phi - \theta_{n-1}))) \quad [13]$$

595 Where γ_m is the magnitude of adaptation, γ_s scales the width of adaptation, and *rect* is the half-
 596 wave rectifying function. The responses of the adapted population thus depend on both the current and
 597 previous stimulus (Figure 3A, *Efficient Encoding*):

$$Resp_A(\theta_n, \theta_{n-1}) = R \gamma_A e^{\kappa \cos(\Phi - \theta_n) - 1} \quad [14]$$

598

599 *Unaware decoder*: We first considered a model in which an adapted orientation-encoding
 600 representation is being decoded by an *unaware* readout mechanism (Figure 3B). The likelihood of each
 601 orientation giving rise to the observed response profile across N neurons was estimated assuming
 602 activity was governed by a Poisson process:

$$P_{unaware}(Resp_A|\theta) = \exp\left(\sum_{i=1}^N \log P_{Poisson}(Resp_A^i(\theta); Resp_N^i(\theta))\right) \quad [15]$$

$$P_{Poisson}(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \quad [16]$$

603 Where $Resp_N^i(\theta)$ is the expected response of the unadapted neuron i to stimulus θ and $P_{Poisson}(k; \lambda)$ is
 604 the probability of observing k spikes given an expected firing rate of λ . The decoded orientation is then
 605 the θ giving rise to the maximum likelihood (MLE).

606 *Aware decoder*: In addition to the unaware decoder, we also evaluated the ability of a decoder that was
 607 aware of the current state of adaptation to explain behavior. The *aware* decoder differs from the
 608 *unaware* decoder in that its assumed activity level for each unit is modulated as a function of stimulus
 609 history:

$$P_{aware}(Resp_A|\theta_n; \theta_{n-1}) = \exp\left(\sum_{i=1}^N \log P_{Poisson}(Resp_A^i(\theta_n, \theta_{n-1}), Resp_A^i(\theta_n, \theta_{n-1}))\right) \quad [17]$$

610 Note that here the rate parameter $k \equiv \lambda \equiv Resp_A$ such that the observed and expected values
 611 perfectly align with the presented orientation. $P_{aware}(Resp_A|\theta_n; \theta_{n-1})$ is dependent on sensory history
 612 and is non-biased.

613 *Over-Aware decoder:* Our final decoding scheme we call the *over-aware decoder*. This model can test
 614 whether serial dependence can be achieved without an explicit stage of Bayesian inference introduced
 615 in the next section. The decoder has an assumed adaptation defined by a unique set of free parameters
 616 γ_{m2} and γ_{s2} which shapes a separate gain adaptation:

$$\gamma_{OA}(\theta_{n-1}) = \gamma_N - rect(\gamma_{m2} \cos^3(\gamma_{s2}(\Phi - \theta_{n-1}))) \quad [18]$$

617 which in turn shapes the response profile of $Resp_{OA}$ in the same manner as $Resp_A$. The likelihood profile
 618 is then defined as:

$$P_{over-aware}(Resp_A|\theta) = \exp\left(\sum_{i=1}^N \log P_{Poisson}(Resp_A^i(\theta); Resp_{OA}^i(\theta, \theta_{n-1}))\right) \quad [19]$$

620 where our expected (assumed) rate λ is designated by $Resp_{OA}$. By having a larger assumed adaptation
 621 than implemented at encoding (through either $\gamma_{m2} > \gamma_m$ or $\gamma_{s2} > \gamma_s$) the net effect of the over-aware
 622 decoder should be behavioral attraction.

623 *Bayesian Inference:* In addition, we explored the effect of applying an explicit Bayesian prior based on
 624 temporal contiguity to the likelihood functions derived from these different readout schemes. This type
 625 of prior has been previously used to explain behavioral biases without considering how encoding might
 626 also be affected by stimulus history³⁵. Specifically, the prior is defined by the transition probability
 627 between consecutive stimuli and is defined as a mixture model of a circular Gaussian and a uniform
 628 distribution:

$$P_T(\theta_n|\theta_{n-1}) = \frac{1}{Z} e^{-\frac{angle(\theta, \theta_{n-1})^2}{2\psi^2}} \quad [20]$$

$$P_{Bayesian}(\theta_n|\theta_{n-1}) = P_{SAME}P_T(\theta|\theta_{n-1}) + (1 - P_{SAME})\frac{1}{2\pi} \quad [21]$$

630 With P_{SAME} set to 0.64 (as found empirically in³⁵), Z as a normalization constant so P_T integrates
 631 to 1, and ψ is a free parameter describing the variance of the transition distribution. This prior (Figure
 632 3C, black line) is multiplied by the *unaware* likelihood (Figure 3C, yellow dashed-line): to get the
 633 posterior estimate of our *Bayesian-unaware* decoder (Figure 3C, yellow solid-line):

$$P_{Bayesian-unaware}(\theta_n|Resp_A; \theta_{n-1}) = P_{Bayesian}(\theta|\theta_{n-1})P_{unaware}(Resp_A|\theta_n) \quad [22]$$

635 We can additionally examine a *Bayesian-aware* decoder by substituting its respective likelihood
636 function. We did not examine a *Bayesian-over-aware* model so that all decoding models would have the
637 same number of free parameters and so that we could directly evaluate the need for an explicit prior.

638 *Model Fitting:* The *encoding* stage of the model has two free parameters and for each subject these
639 parameters were optimized to minimize the residual sum of squares (RSS) between our measured fMRI
640 decoding errors and the decoding errors of our *unaware* decoder. For simplicity we only fit our model to
641 decoding errors from V3 as it had the highest SNR, but other early visual ROIs showed similar results.
642 After fitting the *encoding* stage of the model, we then separately fit the three competing *decoding*
643 models to best account for the behavioral data: *Bayes-unaware*, *Bayes-aware*, and *over-aware* (two free
644 parameters each). The output of this readout stage was treated as the behavioral bias (μ) and the free
645 parameters were optimized to maximize the likelihood of the observed responses (assuming constant
646 standard deviation σ estimated empirically for each participant). For the purposes of fitting the model,
647 the firing rates of the modelled neurons were deterministic (no noise process). Having noiseless activity
648 had no effect on the expected bias (verified with additional simulations) and served to make model
649 fitting more reliable and less computationally intensive. Both stages of the model were fit using the
650 same cross-validation groups as our neural decoder. To ensure all models had a sufficient chance of
651 achieving a good fit to behavioral data, we implemented a grid search sampling 30 values along the
652 range of each variable explored (900 locations total) followed by a local search algorithm (Nelder-Mead)
653 around the most successful grid point. We found dense sampling of the initial parameter space was
654 especially important for our *Bayes-unaware* model.

655 *Model Evaluation:* For bias of neural and behavioral responses, we evaluated the performance of the
656 two stages of our model separately. These stages must be evaluated in a qualitatively different manner
657 as the neural data gives us an orientation estimate for each trial while the behavioral data consists of
658 binary responses. For the encoding stage, we quantified how well the output of our *unaware* decoder
659 predicted the raw errors of our BOLD decoder using circular correlation. The performance of this model
660 was contrasted with the true presented orientation which is analogous to the representation of an
661 unadapted population. We additionally computed the variance of the neural decoding errors explained
662 by the model bias (R^2). For the decoding stage of our model, we compared the log-likelihood of
663 observed responses for each model.

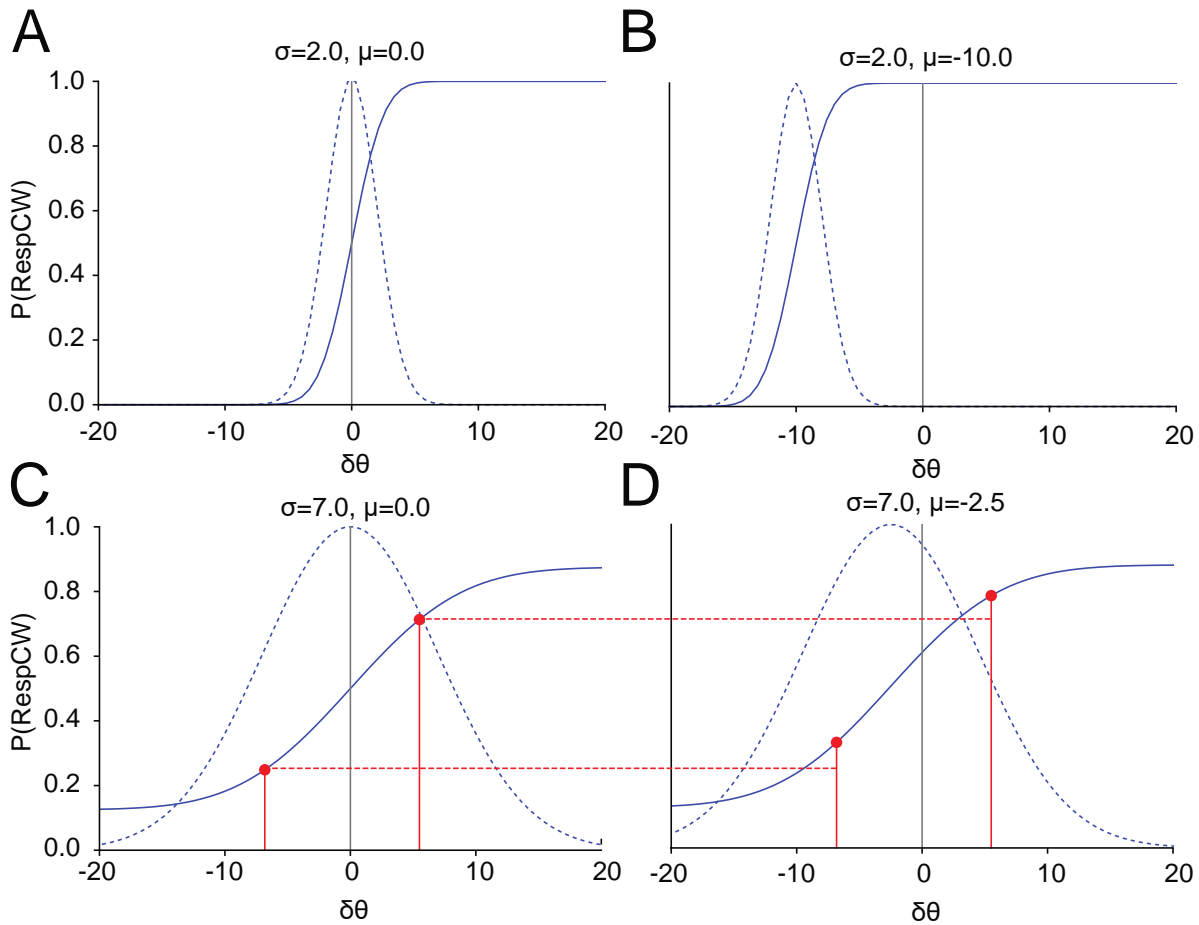
664 We additionally estimated the variance of our models using neurons with rates generated by a
665 Poisson process. The average bias was unaffected by allowing random fluctuations in activity, but the
666 trial-to-trial variance increased. To get a stable estimate, we simulated 1000 trials for each set of
667 parameters estimated for a cross-validation loop for each participant and pooled these outputs. We
668 compared the overall variance of our models to our single parameter estimate of participant precision
669 using Jensen-Shannon divergence. We additionally examined relative precision of our model for close
670 and far trials in the same manner as participant responses and decoding errors ([Response Precision](#)).

671 [Data/Code Availability \(upon acceptance for publication\)](#)

672 Code for processing raw data as well as for analyzing decoded representations can be found
673 here (GITHUB). This includes all processing performed on our BOLD data and our implementation of the
674 Bayesian decoder in Python code for running our model as well as the data used to fit models can be

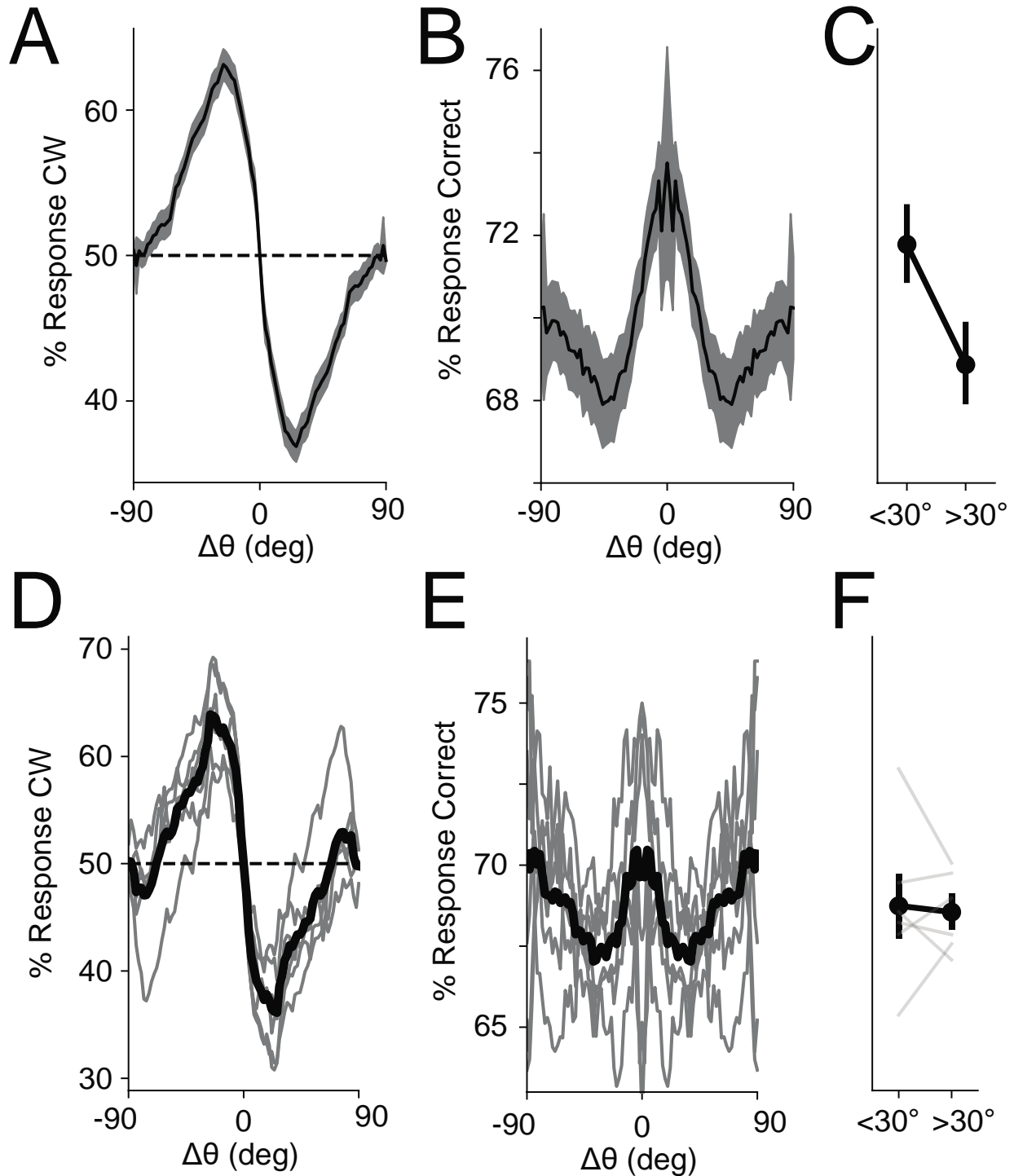
675 found here (GITHUB). Data used in this study will be posted on the first author's Open Science
676 Framework repository (REPOSITORY LINK).

Supplemental Figure 1



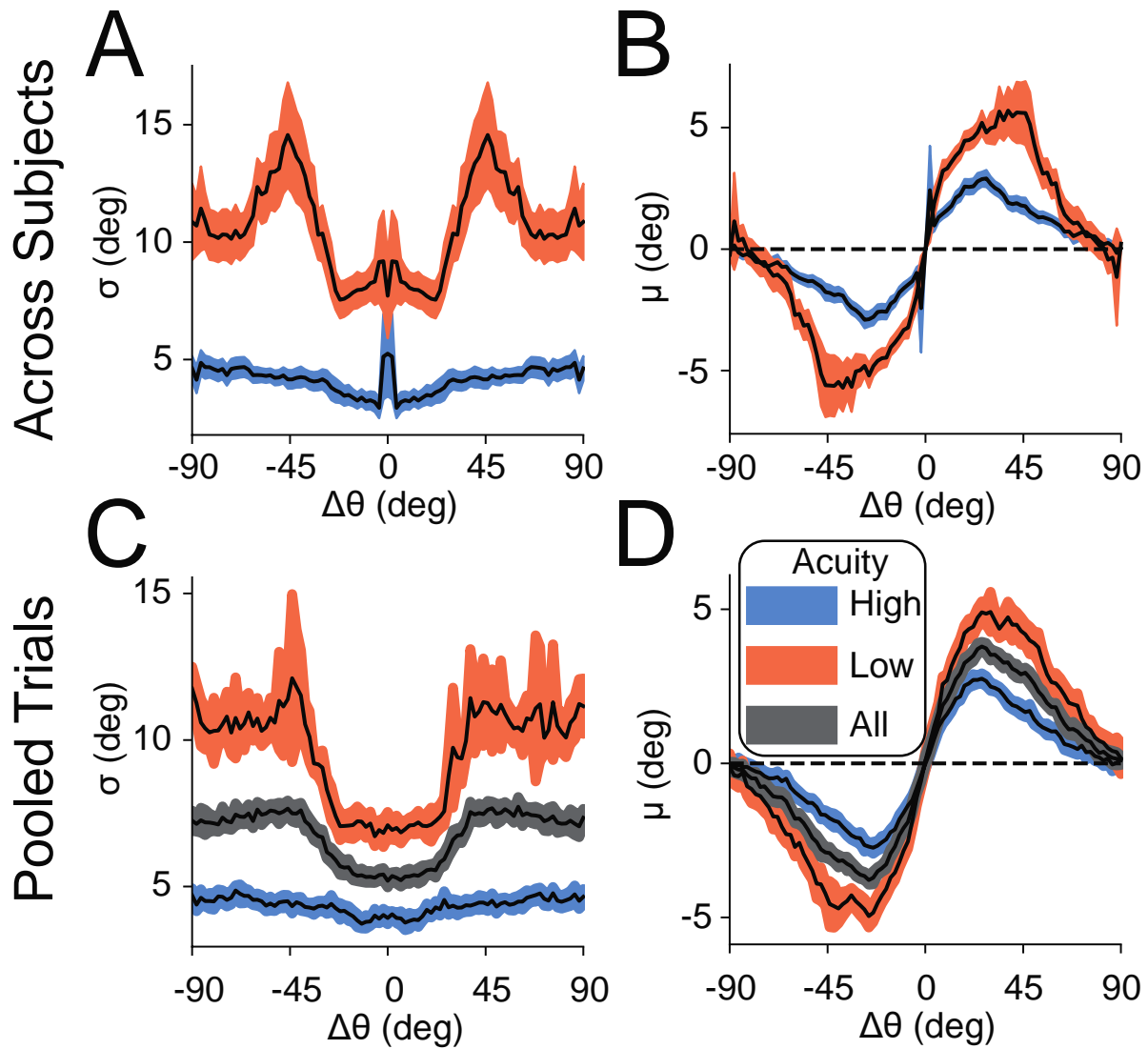
Supplemental Figure 1 Response model. Encoding of stimulus is assumed to be a noisy process whereby the distribution of encoded orientation is described by a Gaussian pdf with mean μ and standard deviation σ . Dashed line is pdf and solid line is the cdf of encoding distribution. Note that participants are reporting the probes orientation relative to the stimulus so more frequent CW responses would correspond to a CCW perceptual bias. **A:** Example estimation curve with no bias and a very small σ . If the difficulty was set to $\delta\theta=6^\circ$ (3 sd) than this participant would get essentially all (99.7%) trials correct. **B:** Estimation curve with a $\mu=-10$, this participant would respond CW on almost every trial. **C-D:** Realistic encoding curves. To aid with fitting and to best describe responses, a constant guess rate of 25% was included in the response model fit to participant responses. **C:** An unbiased distribution with two theoretical stimuli on which the participant responded CW. The left response $\delta\theta=-6^\circ$ is incorrect. **D:** A CCW biased distribution results in a higher likelihood for all CW responses.

Supplemental Figure 2



Supplemental Figure 2 Raw responses. **A-C** behavioral participants. **A**: % response CW as a function of $\Delta\theta$. Note the opposite direction of effect as a CW response means the stimulus was perceived to be CCW of the probe. Shading is SEM across participants **B**: % correct as a function of $\Delta\theta$. Note that as with all analyses, trials without a response are excluded. **C**: % correct following close or far stimuli. Close sequences led to significantly more correct trials $T(45)=3.54$, $p=.0005$. **D-F** fMRI participants. Analysis the same but showing individual participants. No significant difference for accuracy between close and far stimuli $p=0.40$.

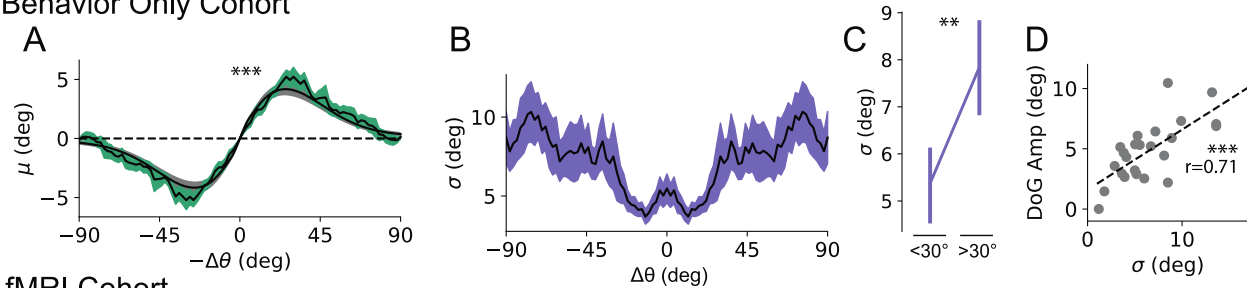
Supplemental Figure 3



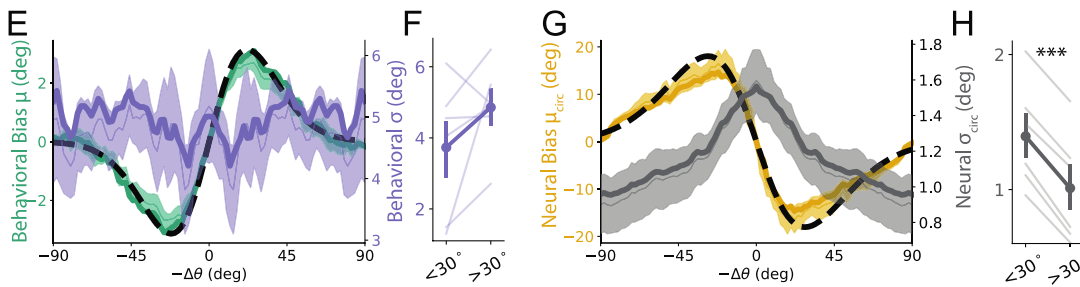
Supplemental Figure 3 Median split bias/variance. **A-B:** average (\pm SEM across participants) across participants. **A:** Model estimated variance for high and low precision participants. **B:** Model estimated bias is larger for less precise ($5.89 \pm 0.52^\circ$) than more precise ($3.57 \pm 0.72^\circ$) participants, $T(44)=2.5$, $p=.007$, unpaired 1-tailed t-test on DoG fits. **C-D:** pooled analysis (\pm 95% bootstrapped CI). **C:** same as A. Insert shows high acuity participants on own axis. **D:** Same as B.

Supplemental Figure 4

Behavior Only Cohort

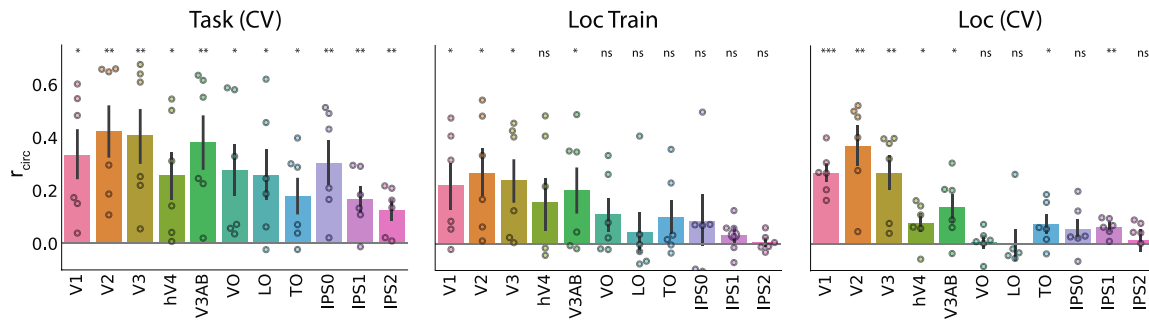


fMRI Cohort



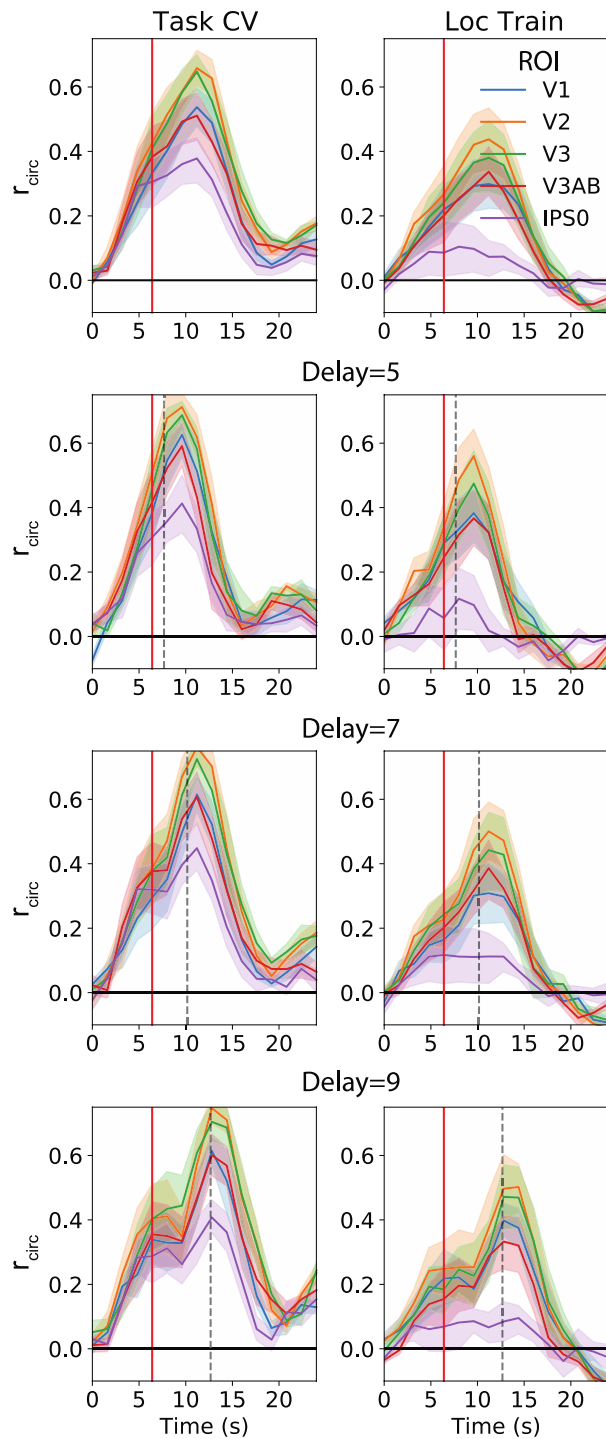
Supplemental Figure 4 Random sequences only. **A-C**: Some participants completed blocks of trials where orientation was not independent across time. Specifically, trials were sorted such that the distance between consecutive stimuli was a multiple of 22.5° (with some jittering). This was intended to maximize our sensitivity to detecting serial dependence as prior experiments in our lab have shown serial dependence typically peaks around 25 degrees. Our behavioral effects ended up being robust and we later opted to just have stimuli be independent across time. Despite no participants overtly noticing any pattern, it is possible that this contrived setup somehow contributed to the behavioral trial history effects that we observed. To assess the impact of this manipulation, we separately analyzed data from only those participants who completed the task with independent stimulus sequences. This cohort ($N=25$) had an average accuracy of $70.46 \pm 1.14^\circ$ at an average $\delta\theta$ of $4.97 \pm 0.35^\circ$. **A**: Serial dependence. The average amplitude when parameterized with a DoG was still significantly greater than 0 (amp= 4.71 ± 0.49 , $T(23) = 9.4$, $p=2.4 \times 10^{-9}$; width 0.027 ± 0.0019 , FWHM $43.68 \pm 1.86^\circ$, (mean \pm SEM). **B-C**: Response variance. Responses were still significantly more precise following similar stimuli ($t(24)=-2.66$, $p=0.01$). **D**: bias and variance were still positively correlated across participants ($r(22)=0.71$, $p=0.00005$). **D-G**: As with our behavior only cohort, some fMRI participants completed blocks of trials where trials were not independent across time. We re-ran a series of control analysis excluding these blocks and found little change to our main findings. **E**: Responses were still systematically attracted to the previous stimulus (DoG Amp: 3.25 ± 0.34 , $T(5)=8.85$, $p=1.53 \times 10^{-4}$; DoG Width: 36.1 ± 2.9). **F**: Response variance was no longer significantly smaller following small changes but was trending in that direction ($T(5)=-1.55$, $p=.09$). **F-G**: Decoded representations showed the same robust pattern of repulsive bias and uncertainty as the full dataset. Together this suggests that our results were not somehow corrupted by the set of trials in which stimuli were not independent.

Supplemental Figure 5



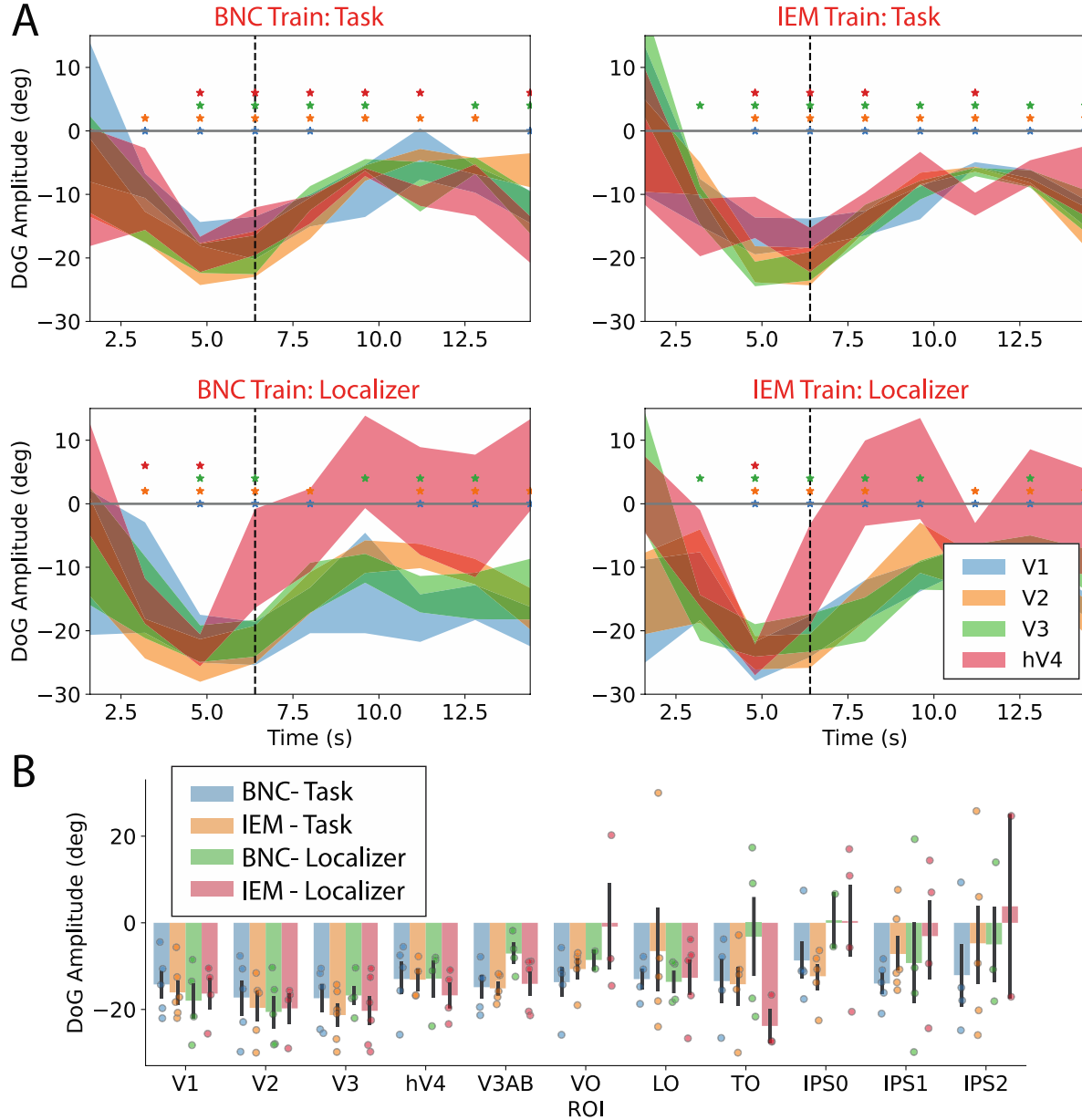
Supplemental Figure 5 Model Performance. Relative performance of model trained and tested on task (**left**), on a separate localizer paradigm and tested on the task (**middle**), or trained and tested within the localizer task (**right**). For the model trained and tested on the localizer data, we could not use orientation information during voxel selection as this would be a circular analysis. Instead, we performed a 75% voxel threshold on donut selectivity for each ROI. See [Voxel Selection](#) for selection process for localizers tested on task data. *, $p < 0.05$; **, $p < 0.01$; ***, $p < 0.001$; t-tests on Fisher transformed r_{circ} .

Supplemental Figure 6



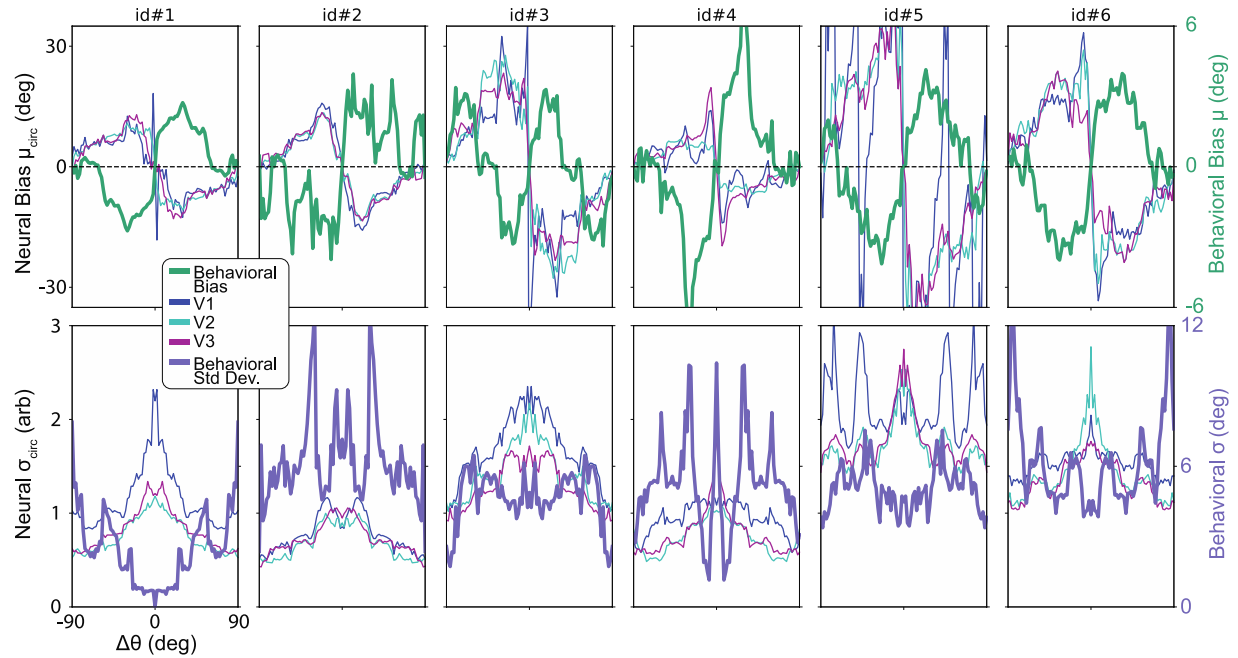
Supplemental Figure 6 Model Performance across time. **Left column:** model trained on task. **Right column:** model trained on localizer and tested on all task TRs. **Top Row:** all trials. **Rows 2-4:** subset of trials corresponding to delays of 5, 7 and 9s. Shaded lines depict average r_{circ} across participants (\pm SEM across participants) for 5 ROIs (see legend). Dashed vertical line is average delay time for a given group. Red vertical line is central TR used in main analyses. Time is not shifted to account for hemodynamic lag so even the probe on the shortest delay trials should not affect signal measured at red line.

Supplemental Figure 7



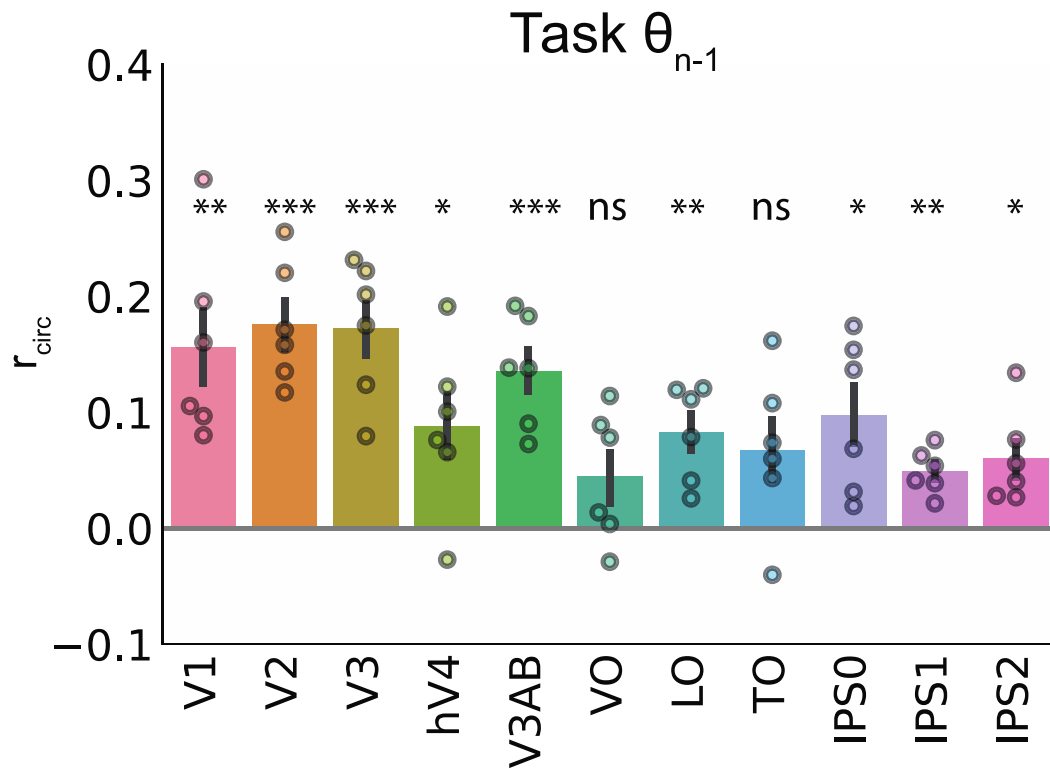
Supplemental Figure 7: Decoded bias across time and ROI. **A:** Decoded bias is significantly repulsive (and never attractive) across the extent of the trial when parameterized with a DoG. For completeness we show two different decoders (“BNC”, [Bayesian Noise Correlation](#); and “IEM”, [Inverted Encoding Model](#)) trained on both the task and localizer data. Time points represent middle of sliding 4 TR window. *, $p < .01$, uncorrected. **B:** Decoded bias is generally repulsive (and never attractive) across all ROIs and decoding techniques for TR window centered at 6.4s (indicated with dashed line in **A**).

Supplemental Figure 8



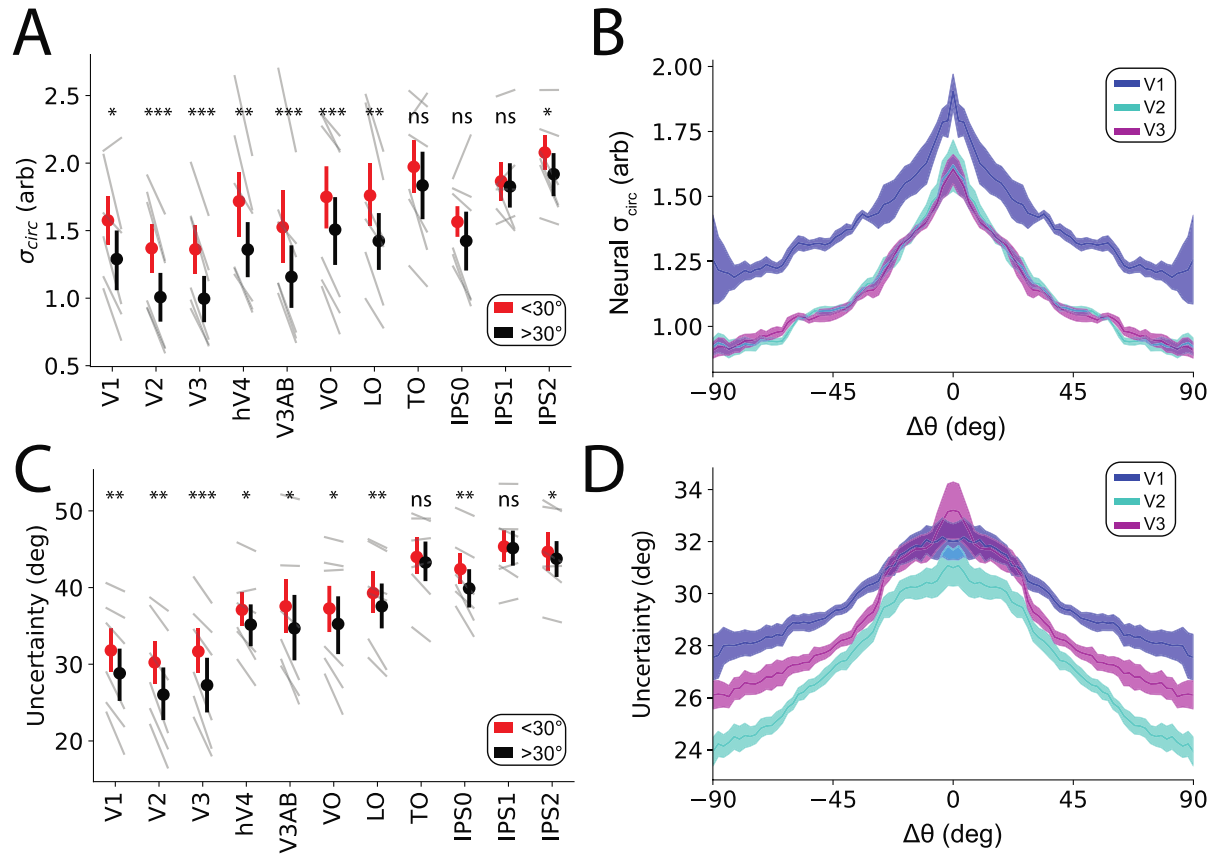
Supplemental Figure 8 Decoded and behavioral bias (**top**) and variance (**bottom**) for individual participants. **Left axis:** Neural data for ROIs V1, V2, and V3 (see **legend**). Decoded orientation is clearly repelled in all participants in V1-V3 and neural σ generally peaks at $\Delta\theta=0$. **Right axes:** Behavioral data. Responses are clearly attracted for all participants. Note how participant id#3 has peripheral *repulsion* from very distant stimuli.

Supplemental Figure 9



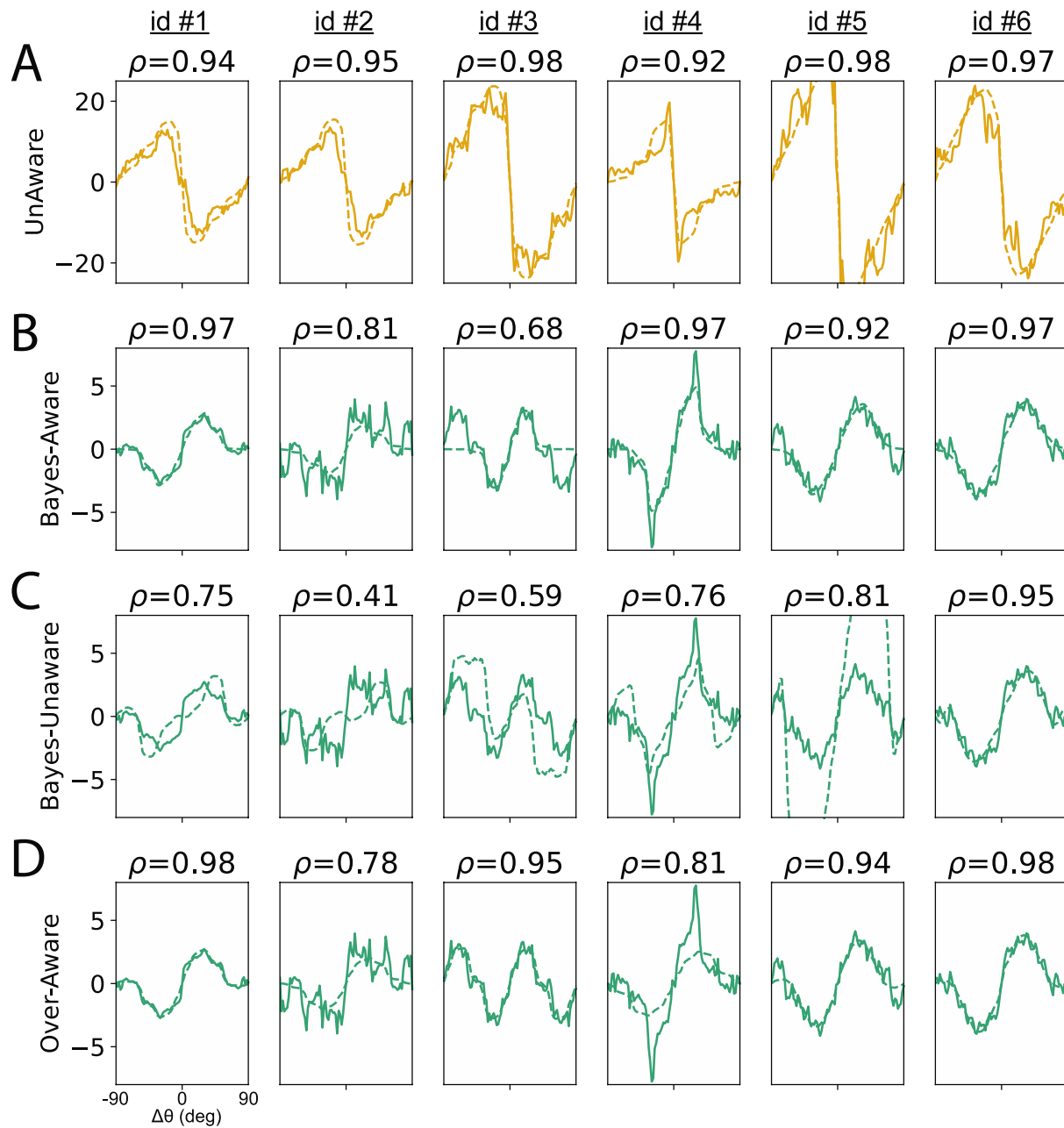
Supplemental Figure 9 Decoding performance for model trained and tested on task data to decode previous trial's stimulus (θ_{n-1}). Performance was significantly above chance in most ROIs. ns, not significant, *, $p < .05$, **, $p < .01$, ***, $p < .001$.

Supplemental Figure 10



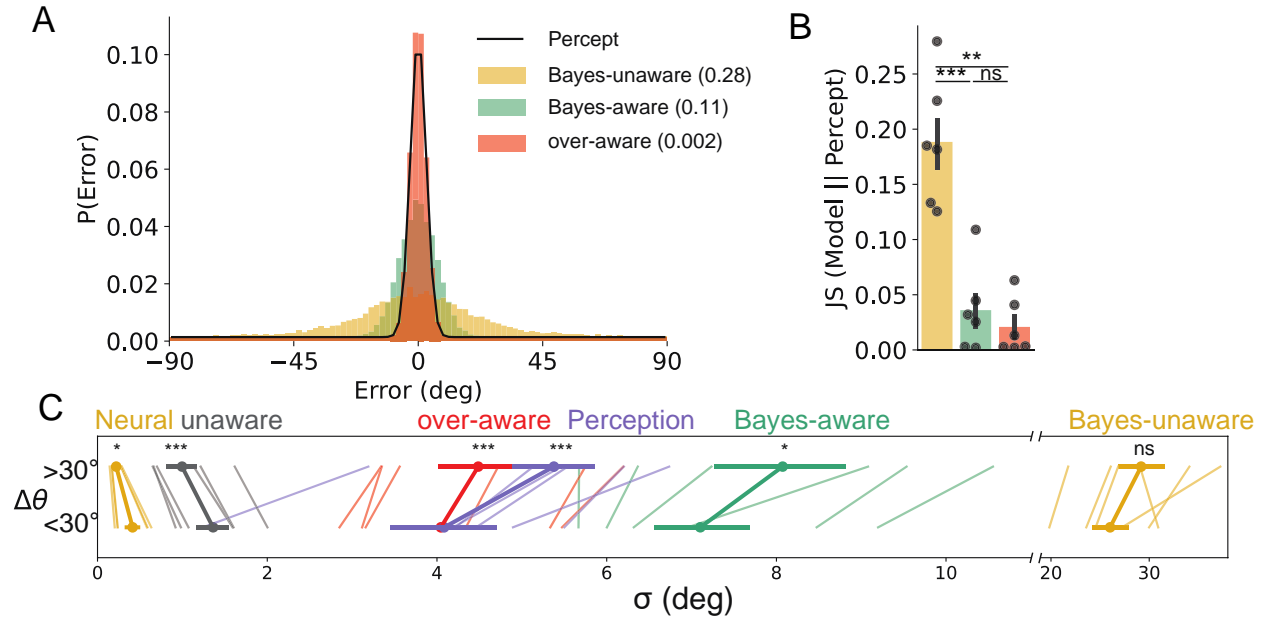
Supplemental Figure 10 Decoded uncertainty as a function of $\Delta\theta$ across ROIs. **A:** σ_{circ} of decoding errors is significantly greater for close ($<30^\circ$) versus far ($>30^\circ$) stimuli across early visual ROIs. Points and error bars are mean \pm SEM across participants; gray lines depict individual participants. **B:** Sliding σ_{circ} for V1-V3 shows a monotonic relationship. **C-D:** Same as **A-B** but measuring uncertainty directly measured from the single trial posterior (see [Neural Variance](#)). Results are qualitatively very similar for both techniques. ns, not significant, *, $p < .05$, **, $p < .01$, ***, $p < .001$.

Supplemental Figure 11



Supplemental Figure 11 Model fits for individual participants (same order as Figure S8). Solid lines correspond to empirical neural (yellow) or behavioral (green) bias; dashed lines correspond to model fits to BOLD decoding bias (Unaware model, A) or behavior (B-D). Model fits plotted are average of noiseless biases generated by models fit to each CV fold. Note that a models are fit to raw data, not binned data presented here. Pearson's correlations are reported above each fit between binned and model estimated bias.

Supplemental Figure 12



Supplemental Figure 12 Model Performance. **A:** Distribution of empirically predicted response errors (black line) and simulated model fits for an example participant along with associated Jensen-Shannon divergences. **B:** The *Bayes-unaware* model provided a significantly worse fit to empirical uncertainty than either “aware” model when assessed across participants. **C:** Visualization of all uncertainties split as a function of close and far stimuli. The *unaware* model was significantly less precise following small changes matching neural decoding. The two “aware” models were significantly more precise following small changes matching perception. The *Bayes-unaware* model did not have significant modulation of decoding uncertainty and had an average uncertainty that was on average 5x that of perception. ns, not significant, *, $p < .05$, **, $p < .01$, ***, $p < .001$.

Table 1

Average fit coefficients for models averaged across CV fits (\pm SEM across participants) shown in bold. Other parameters either fixed values, drawn from fit to encoding stage for given participant, or are not utilized for a particular model (N/A).

	Fit To:	BOLD Decoder	Behavior		
Stage:		Unaware	Bayes unaware (Prior*unaware)	Bayes aware (Prior*aware)	Over-aware
Encoding	γ_m	0.81 \pm 0.04			
	γ_s	0.56 \pm 0.16			
Decoding	γ_{m2}	0	0	γ_m	0.69 \pm 0.09
	γ_{s2}	1	1	γ_s	0.68 \pm 0.11
Bayes	R	5	1.95 \pm 0.57	0.17 \pm 0.03	5
	ψ	N/A	0.60 \pm 0.05	0.86 \pm 0.05	N/A