1 **Supplementary Materials**

2 **Attention to colors induces surround suppression at category boundaries**

3 **Ming W.H. Fang, Mark W. Becker, Taosheng Liu[*]**

4

5 **Neural model for surround suppression at categorical boundaries**

6

7 Our goal here is to build the simplest model that is informed by physiological data to produce

8 the attentional profile measured in our psychophysical experiments. The advantage of this

9 approach, compared to a full-blown model such as multi-layered neural network model, is that

10 we have a much better understanding of how model parameters impacts its behavior. However,

11 simplicity is only achievable by ignoring many physiological details and as such, our model is

12 more of a proof-of-concept than a complete description of the physiological processes.

13 Nevertheless, such a model can still give useful insights regarding the neural mechanisms of

14 attention.

15

16 The model contains a bank of identical, uniformly distributed, color-tuned neurons spanning

17 the 360° space defined by the color wheel. Each neuron's tuning curve is assumed to be a

18 circular Gaussian function (von Mises function)

19
$$R_{ij} = \frac{e^{\kappa \cos(\theta_j - \mu_i)}}{2\pi \, I_0(\kappa)} \cdot (A - s) \; + s \qquad (1),$$

20 where $R_{ij}$ is the i-th neuron's response to a colored dot $\theta_j$, $\kappa$ is the concentration parameter that

21 controls the spread of the tuning function, and $\mu_i$ is the neuron's preferred color. $I_0(\kappa)$ is the

22 Bessel function of order 0. Parameter $A$ denotes the firing rate to the preferred color, and $s$

23 represents the spontaneous firing rate. The model contains 90 neurons with $\kappa$ =12, $s$ =10

24 spikes/s, $A$=40 spikes/s. The parameter values are based on relevant physiological findings

25 (see Table S1 for the full list of parameters and their values). For simplicity, we assumed no

26 neural noise or inter-neuronal correlation.

27

28 In a simulated trial of the 2-IFC task, the model is "shown" two random dot color stimuli, a

29 noise pattern with random colors and a signal pattern at a particular color coherence. Each dot

30   independently evokes responses across all neurons, which are computed using Eq. (1). Each
31   neuron's response is determined by averaging its responses to all dots in the stimulus, which
32   is computed by

$$R_i = \frac{1}{N} \sum_{j=1}^{N} R_{ij} \qquad (2),$$

34   where $N$ is the total number of color dots in each stimuli array (fixed at 100), and $R_i$ is the
35   neuron's average response to all dots. The response across all neurons to a dot array thus
36   constitutes a population neural response and is the basis of the model's decision. For each
37   stimulus interval, we fitted a Gaussian template (Eq. (1)) to the population response using
38   maximum likelihood estimation (MLE), which had four free parameters – amplitude ($A$), mean
39   ($\mu_i$), variance ($\kappa$), and intercept ($s$). As the 2-IFC task requires participants to detect a stronger
40   color signal, we used the estimated amplitude ($A$) as the decision criterion. The model simply
41   chose the stimulus interval with a higher amplitude estimate as the target.

42

43   We first performed baseline simulations by presenting the model with color stimuli of different
44   coherence levels. The model's choice for each trial was recorded and the proportion correct
45   rate was calculated. For all results presented here, we simulated 2000 trials for each condition.
46   In this baseline (neutral) condition, the model performed better with higher color coherence,
47   similar to human observers (Fig. S1a). We also checked the population response for stimuli of
48   different coherence levels and found it to increase monotonically with coherence (Fig. S1b).
49   This increase in population response thus reflected the increase in the signal strength and was
50   appropriately registered by the model. For the main simulations, we fixed the color coherence
51   at 0.1 as it produces an intermediate performance level in the neutral condition. Indeed, this
52   coherence level was comparable to coherence thresholds measured in our human participants
53   (cf. Fig 5).

54

55   For attention condition, we first simulated the experiment with a pure feature-similarity gain
56   modulation, which was implemented as a linear function:

$$G_{FSG} = b - a * |(\theta_{attend} - \theta_{target})| \qquad (3),$$

58    where $G_{FSG}$ is the gain factor, *b* denotes the intercept of the attentional gain, and *a* represents

59    the slope. This equation expresses the FSG principle: the attentional gain factor for a target

60    feature (i.e., $\theta_{target}$) depends on its difference (similarity) to the attended feature (i.e., $\theta_{attend}$).

61    Without losing generality, we assumed $\theta_{attend}$ =0°. Values of *a* and *b* were based on published

62    values from monkey MT (Martinez-Trujillo & Treue, 2004, see Table S1). The FSG modulation

63    led to a simple scaling of all the tuning curves (Fig. S2a). To facilitate the understanding of the

64    model behavior, we plotted the model's population response to a color signal under the

65    attention and neutral condition for a few selected cue-target offset (Fig. S2b). As can be seen,

66    compared to the neutral condition, population response for the 0° target was higher and

67    gradually declined as target deviated more from 0° such that at large offsets, it became lower

68    than the neutral condition. This monotonic decline of population response underlies the

69    model's monotonic cueing effect (Fig. 8d).

70

71    Next, we implemented a hybrid model by combining the FSG modulation with neuronal tuning

72    shifts (Fig. S3a). The FSG factor is calculated using Eq. 3 above. The magnitude of neurons'

73    tuning shift towards the cued color, *M*, is calculated by a piece-wise linear function,

74
$$
M = \begin{cases} 0.5 \cdot (\theta_i - \theta_{attend}), & if \ |\theta_i| \leq w, \\ 3 \cdot w \cdot sgn\,(\theta_i) - 2.5 \cdot \theta_i, & if \ w < |\theta_i| \leq 1.2w, \\ 0, & if \ 1.2w < |\theta_i| \leq 180° \end{cases} \tag{4}
$$

75    where *w* denotes the boundary (40° in current case), $\theta_i$ is the neuron's original tuning

76    preference, and $\theta_{attend}$ denotes the attended color (fixed at 0°), and *sgn* is the sign function.

77    This results in a larger shift as neurons move further away from the attended feature followed

78    by a reduced shift beyond the category boundary (see Fig. 8f). Once *M* declines to 0, the

79    tuning shift would stop. Under this scenario, neuronal responses were calculated in the same

80    fashion as in Eq. (1), except that neuron's preferred color ($\mu_i$ ), was replaced by ($\mu_i$-*M*),

81    representing a shift in tuning preference.

82

83    The population responses under this hybrid modulation exhibited a non-monotonic profile.

84    Critically, there was a suppression of population response for the boundary color compared to

85    neutral baseline (Fig. S3b, 40°), which was not seen in the FSG only condition (cf. Fig. S2b,

86  40°). This was followed by a relative increase in population response at 60°, signifying a
87  rebound. Finally, for large feature offsets such as 140°, there was a further suppression, as a
88  result of FSG modulation. These qualitative observations on the population responses were
89  registered by our model using the simple read-out rule described earlier, resulting in a hybrid of
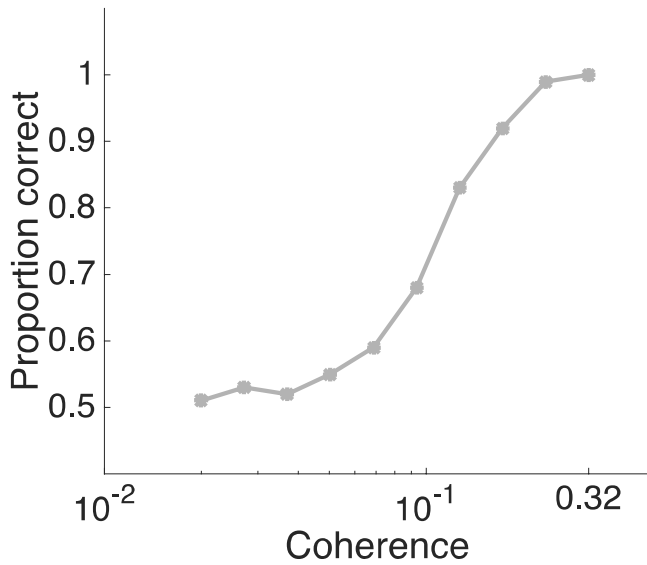90  surround suppression and feature-similarity gain modulation in its performance (see Fig. 8g).
91
92  To verify whether surround suppression can appear at the categorical boundaries, we
93  simulated the experiment with a number of different category width (e.g., ±30°, ±50° ±60°, ±70°,
94  ±80° boundaries) and observed the suppressive surround occurring at the category boundaries.
95  We also explored different shifting parameters that control the exact shape of the shifting
96  function (Fig. 8f) and found that as long as the tuning shift returns to zero beyond the category
97  boundary relative quickly, the model produces a surround suppression at the boundaries. For
98  example, the slope of the declining portion of tuning shift beyond the boundary can be
99  shallower. We also used a sinusoidal tuning shifting function and found similar results with the
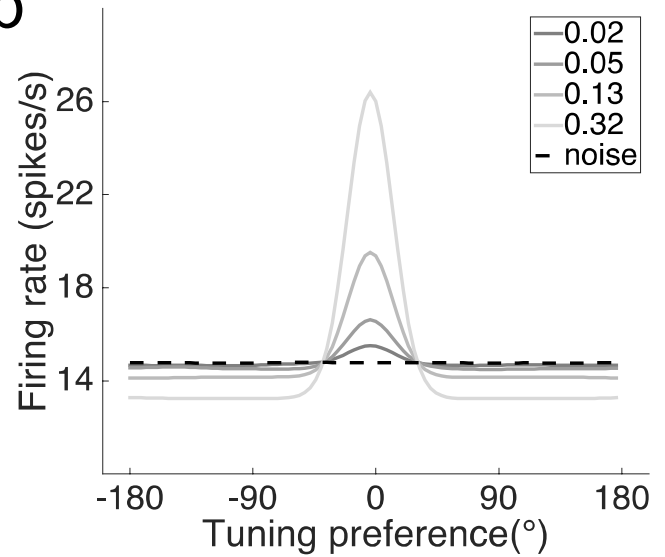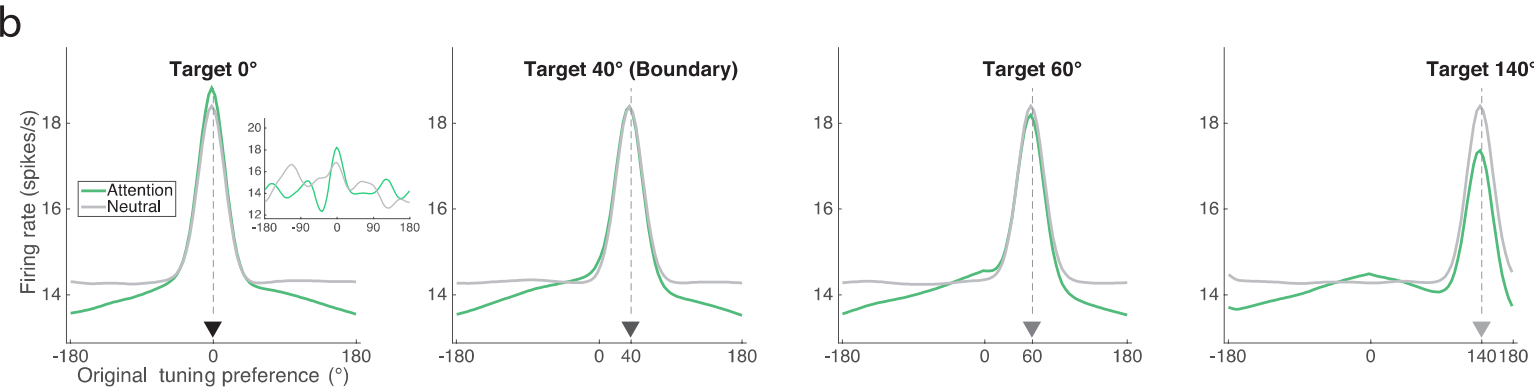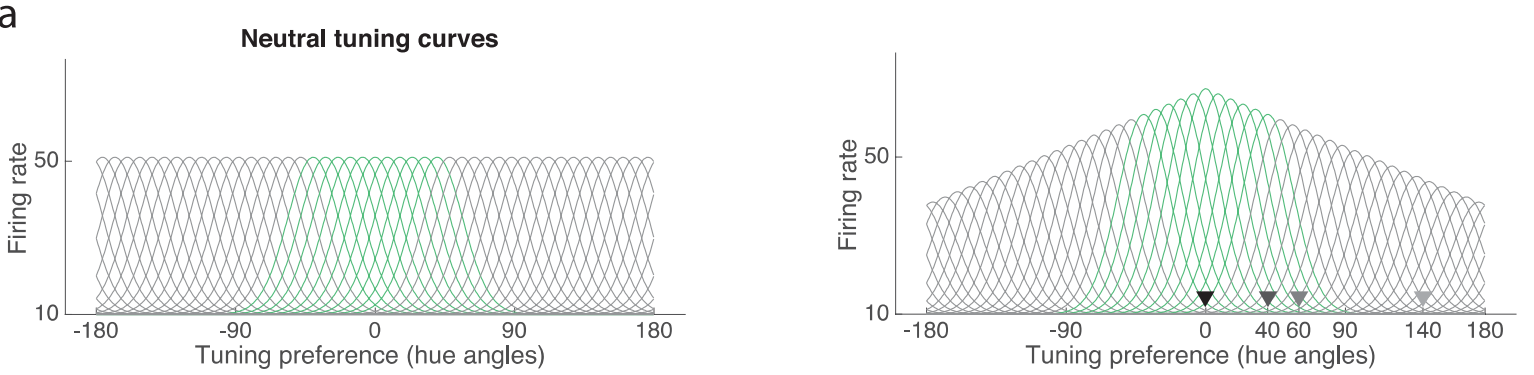100 piece-wise linear function in Eq. 4.
101
102

a



b



Figure S1. The model's neutral (baseline) performance and population responses under variable coherence levels. a). The model performed better as the coherence of color stimuli increased. b) Average population response across trials for a few selected coherence levels (gray lines). The dashed black line denotes the average population response for the noise stimuli.

106
107
108
109
110
111
112
113
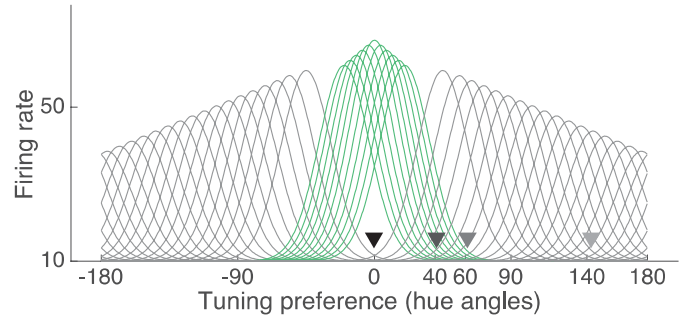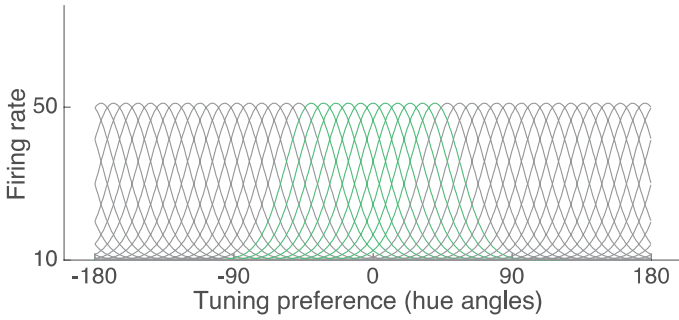114
115
116
117
118
119
120
121
122
123

124

125

126

127

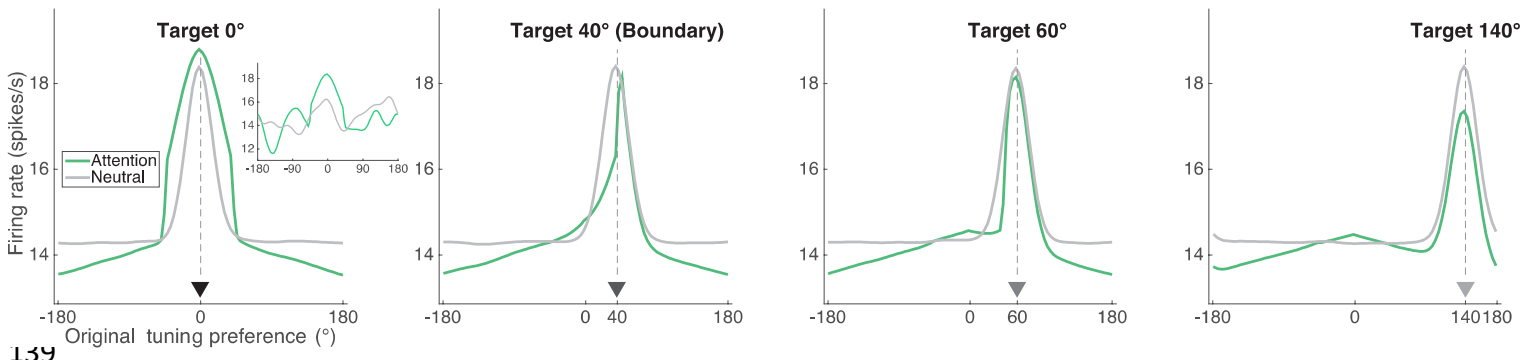128

a

**Neutral tuning curves**



b



129

130  Figure S2. Pure FSG modulation on neural responses. **a)**. Groups of neuronal tuning curves in

131  the neutral condition (left panel) and the attention condition due to FSG modulation (right

132  panel). Attended feature is assumed to be at 0°. **b)**. Example average population responses to

133  a color signal at four cue-target offsets. There was a monotonic decrease in the population

134  responses (green curves) compared to neutral baselines (gray curves). The population

135  responses were averaged across 2000 trials and thus appear quite smooth. Population

136  responses on individual trials were much noisier. The inset in the left most panel shows

137  population responses on a single trial.

138

Figure S3. Combined modulation of FSG and tuning shift on neural responses. **a)**. In addition to FSG modulation, neurons also shifted toward the category center. **b)**. Example population responses to a color signal at four cue-target offset. Note the population responses in the attention condition (green curves) changed non-monotonically compared to neutral baselines (gray curves). The smooth population responses were averaged across 2000 trials. Examples of population responses on a single trial are shown in the inset of the left most panel.

156

157    Table 1. Parameters and their values used in the neural model simulation.

| Name | Value | Description |
|------|-------|-------------|
| κ | 12 | Single neuron's tuning bandwidth, equivalent to a bandwidth of ~39°, similar to previously reported value (Conway, Moeller, & Tsao, 2007) |
| b | 1.0372 | Intercept of attentional gain factor in the FSG model, based on values reported by Martinez-Trujillo, & Treue (2004) |
| a | 0.00093 | Slope of attentional gain factor in the FSG model, based on values reported by Martinez-Trujillo, & Treue (2004) |
| N | 100 | Number of colored dots in the simulation |
| s | 10 spikes/s | Neuron's spontaneous firing rate |
| A | 40 spikes/s | Neuron's firing rate to its preferred color |

158

159

160