

# Opponent surrounds explain diversity of contextual phenomena across visual modalities

David A. Mély<sup>1,2,3</sup> & Drew Linsley<sup>1,2,3</sup> & Thomas Serre<sup>1,2,3</sup>

<sup>1</sup>*Department of Cognitive, Linguistic & Psychological Sciences*

<sup>2</sup>*Brown Institute for Brain Science*

<sup>3</sup>*Brown University, Providence, RI 02912, USA.*

**Running head:** Opponent surrounds explain diversity of contextual phenomena

**Corresponding Author:** Thomas Serre, Department of Cognitive, Linguistic & Psychological Sciences, Brown University, 190 Thayer St, Providence, RI 02912, USA.

## 1 **Abstract**

2 Context is known to affect how a stimulus is perceived. A variety of illusions have been attributed  
3 to contextual processing — from orientation tilt effects to chromatic induction phenomena, but  
4 their neural underpinnings remain poorly understood. Here, we present a recurrent network model  
5 of classical and extra-classical receptive fields that is constrained by the anatomy and physiology  
6 of the visual cortex. A key feature of the model is the postulated existence of near- vs. far-  
7 extra-classical regions with complementary facilitatory and suppressive contributions to the classical  
8 receptive field. The model accounts for a variety of contextual illusions, reveals commonalities  
9 between seemingly disparate phenomena, and helps organize them into a novel taxonomy. It  
10 explains how center-surround interactions may shift from attraction to repulsion in tilt effects,  
11 and from contrast to assimilation in induction phenomena. The model further explains enhanced  
12 perceptual shifts generated by a class of patterned background stimuli that activate the two opponent  
13 extra-classical regions cooperatively. Overall, the ability of the model to account for the variety and  
14 complexity of contextual illusions provides computational evidence for a novel canonical circuit  
15 that is shared across visual modalities.

16 **Keywords:** extra-classical receptive field, visual cortex, illusion, induction, assimilation, tilt effect.

17 Spatial context has been known to affect perception since at least Aristotle (Eagleman, 2001).  
18 The past several decades of work in visual psychophysics have revealed a plethora of seemingly  
19 disparate contextual phenomena (Series et al., 2003) whereby subtle differences in experimental  
20 conditions yield a wide variety of effects (Figure 1). In the classical tilt illusion (O’Toole and  
21 Wenderoth, 1977; Goddard et al., 2008), the perceived orientation of a center stimulus tilts either  
22 towards or away from that of a surround stimulus, depending on their relative orientations. Many  
23 variants have been tested with a variety of stimulus parameters including spatial frequency, color,  
24 luminance, contrast differences between center and surround stimuli as well as their spatial and  
25 temporal separation (see Clifford, 2014; for a review).

26 Similar effects have been reported in the motion domain – for both direction and speed (Marshak  
27 and Sekuler, 1979; Murakami and Shimojo, 1993; 1996; Kim and Wilson, 1997). In color induction,  
28 both the spatial frequency and phase of the surround controls the direction of the perceived shift in  
29 hue of a center stimulus relative to that of the surround (Smith et al., 2001; Monnier and Shevell,  
30 2003; Shevell and Monnier, 2005). In the disparity domain, a center stimulus appears closer or  
31 further away from an observer, depending on the relative depth and spacing between center and  
32 surround stimuli (Westheimer, 1986; Westheimer and Levi, 1987). While much is known about the  
33 psychological basis of these phenomena, our understanding of the underlying neural mechanisms  
34 remains, at best, fragmentary.

35 A widely held assumption is that such contextual phenomena are mediated in the cortex by extra-classical  
36 receptive field (eCRF) mechanisms (reviewed in Series et al., 2003; Angelucci and Shushruth,  
37 2013): The presentation of a stimulus in the eCRF alone does not typically elicit any response

38 from the neuron but modulates its response to a stimulus presented in the classical receptive field  
39 (CRF). Such center-surround interactions have been reported across visual modalities including  
40 orientation and spatial frequency (DeAngelis et al., 1994), motion (Li et al., 1999; Jones et al.,  
41 2001), color (Schein and Desimone, 1990; Wachtler et al., 2003) and disparity (Bradley and  
42 Andersen, 1998).

43 Although several eCRF models have been developed to describe specific phenomena (reviewed in  
44 Series et al., 2003; Schwartz et al., 2007; Angelucci and Shushruth, 2013; see also Discussion), a  
45 unifying theory, which would integrate disparate aspects of contextual integration and, ultimately,  
46 link primate neurophysiology to human behavior, is still lacking. We have thus developed a  
47 large-scale recurrent network model of classical and extra-classical receptive fields that distinguishes  
48 itself from previous work – allowing us to simulate realistic cortical responses to a variety of  
49 full-field, real-world, contextual stimuli defined across visual modalities (we model orientation,  
50 color, motion, and binocular disparity). The model is constrained by anatomical data and shown  
51 in our experiments to be consistent with V1 neurophysiology. A key feature of the model is the  
52 postulated existence of near vs. far extra-classical eCRF regions with complementary contributions  
53 (facilitatory vs. suppressive) to the CRF response. Using an ideal neural observer, we show that  
54 the model is consistent with human behavioral responses for a variety of contextual phenomena –  
55 revealing commonalities between seemingly disparate phenomena and helping to establish a novel  
56 taxonomy of contextual illusions.

## 57 **Results**

58 The visual cortex is modeled as a dense, regular topographic grid of cortical (hyper)columns which  
59 tile the visual field (Figure 2A). Each hypercolumn contains a complete set of units with coinciding  
60 CRFs. Their tuning curves are idealized (see Materials and Methods), and centered at regular  
61 intervals (e.g., between  $0^\circ$  and  $180^\circ$  for orientation-tuned units). For simplicity, we do not take into  
62 account cortical magnification and assume a fixed sampling of the visual field at all eccentricities.  
63 The model takes into account connections both within and across hypercolumns in order to explain  
64 several CRF and eCRF properties. The resulting circuit motif is replicated for every hypercolumn.

### 65 **Intra-columnar recurrent circuits**

66 Recurrent connections (Figure 2A, red connections) within a column (i.e., originating from within  
67 the CRF) include both local excitatory and inhibitory connections. Inhibitory CRF contributions  
68 constitute one of the key mechanisms in an influential model of gain control (divisive normalization,  
69 reviewed in Carandini and Heeger, 2012). This model accounts for cross-orientation normalization  
70 phenomena (when a grating stimulus is masked by another one at any orientation, see Heeger,  
71 1993; Carandini and Heeger, 1994) and was later extended to capture neural population responses,  
72 in order to account for the competitive interactions within a single hypercolumn (Busse et al.,  
73 2009; Sit et al., 2009). A recent optogenetic study demonstrated that the underlying circuits  
74 are recurrent rather than feedforward (Nassi et al., 2015). We have also confirmed that this  
75 form of recurrent intra-columnar inhibition is critical for the model to reproduce these types of  
76 competitive interactions within the CRF and for model units to exhibit a realistic contrast response

77 (see Supplementary Experiments; Figures S1-S5.)

78 Because this form of suppression does not seem to depend on the orientation of the afferent and  
79 target cells, it is often called “untuned” inhibition (Rust et al., 2006). In our model, we speculate  
80 that such untuned inhibitory recurrent connections within hypercolumns exist for all other visual  
81 domains (including color, motion as well as binocular disparity, see Discussion).

82 In addition to short-range inhibitory connections within hypercolumns, the model also incorporates  
83 short-range excitatory connections. In the cortex, such excitatory connections may drive neurons  
84 up to ten times more strongly than their feed-forward inputs (Douglas et al., 1995; Stepanyants  
85 et al., 2008). As suggested by Shushruth et al. (2012), we have found that recurrent excitation  
86 within the CRF is essential to account for some of the more complex aspects of surround suppression  
87 (see Supplementary Experiments, Figures S3 and S5), by placing the column in a regime dominated  
88 by recurrent as opposed to feed-forward inputs. However, experimental data on the selectivity  
89 of these recurrent excitatory connections are scarce. Here, we assume that the corresponding  
90 local excitatory connections within a hypercolumn are only weakly tuned, as perfectly untuned  
91 excitation would effectively “flatten out” population response curves.

## 92 **Inter-columnar recurrent circuits**

93 Expanding the optimal stimulus of a cortical neuron immediately beyond its CRF (also commonly  
94 referred to as the “minimum response field” or mRF) may facilitate its response (Bringuier et al.,  
95 1999; Sengpiel et al., 1997; Sceniak et al., 1999; Angelucci et al., 2002a;b; Briggs and Usrey,

96 2011). The area which covers the CRF and its immediate eCRF is sometimes referred to as “peak  
97 spatial summation area”. It is considered distinct from the CRF because a direct stimulation of  
98 this eCRF region in isolation does not elicit any action potential. Since this region is located  
99 immediately beyond the CRF, we deem it the *near eCRF* (or near surround; green annulus and  
100 connections in Figure 2A).

101 A potential neural substrate for the near eCRF includes the short-range, tuned excitatory networks (Lee  
102 et al., 2016) which span a spatial extent consistent with that of the eCRF facilitation (Angelucci  
103 et al., 2002a;b) and amplify co-occurring local inputs at similar orientations (Sengpiel, 1997;  
104 Sceniak et al., 1999; Angelucci et al., 2002a;b; Briggs and Usrey, 2011). In the model, we assume  
105 that all excitatory connections from other hypercolumns centered in the near surround are tuned,  
106 irrespective of the visual modality (i.e., the stimulus with the preferred orientation, or direction of  
107 motion, etc. in the CRF is also most effective in the near eCRF). Also note that our definition of  
108 the near eCRF is purely anatomical and might thus differ from that of others (e.g., Angelucci and  
109 Shushruth, 2013), whose definition is functional in nature.

110 Expanding the optimal stimulus beyond the near eCRF results in neural suppression (first reported  
111 by Hubel and Wiesel, 1968; as hypercomplex tuning). Critically, the presentation of the suppressing  
112 stimulus in the eCRF alone does not elicit any activity from the recorded cell (see Angelucci  
113 and Shushruth, 2013; for review). The tuned nature of these suppressive mechanisms is well  
114 documented across visual modalities: from orientation (Hubel and Wiesel, 1968; DeAngelis et al.,  
115 1994; Weliky et al., 1995; Petrov et al., 2005; Ozeki et al., 2009) to color (Schein and Desimone,  
116 1990; Wachtler et al., 2003), spatial frequency (DeAngelis et al., 1994), temporal frequency (Li

117 et al., 1999; Jones et al., 2001), motion direction and speed (Allman et al., 1985) as well as  
118 binocular disparity (Bradley and Andersen, 1998).

119 Thus, we also define an *inhibitory far eCRF* (or far surround; blue annulus and connections in  
120 Figure 2A), located immediately beyond the excitatory near eCRF. In our model, a hypercolumn  
121 receives tuned inhibition from hypercolumns centered in its far surround.

122 To summarize, contributions from the eCRF as a whole arise from surround regions with opposite  
123 polarities. We do not assume any gap between the CRF and the near eCRF, nor between the near  
124 eCRF and the far eCRF. The model assumes that the near excitatory and far inhibitory eCRFs do  
125 not overlap spatially based on partial anatomical evidence (Angelucci et al., 2002b). In practice,  
126 we found that allowing these two eCRF regions to overlap did not affect the model's ability to fit  
127 experimental data (see Supplementary Experiments).

128 The first key assumption of the model is that, unlike local recurrent interactions within a hypercolumn,  
129 interactions across hypercolumns are “tuned” as only units that share the same preferred stimulus  
130 are directly connected. The second key assumption of the model is an asymmetry between excitation  
131 and inhibition: In the model, excitation only depends on pre-synaptic activity and is purely additive.  
132 Inhibition, on the other hand, from either the CRF or eCRF, depends on both pre- and post-synaptic  
133 activity, and ultimately results in a combination of subtractive and divisive effects (Carandini and  
134 Heeger, 2012). Similar forms of inhibition have been used in previous recurrent network models  
135 to achieve divisive normalization (Grossberg and Todorović, 1988). In practice, this means that,  
136 given a fixed amount of pre-synaptic inhibition, weakly active units receive less effective inhibition

137 than more active ones. In contrast, any given amount of pre-synaptic excitation results in the same  
138 amount of effective post-synaptic excitation.

### 139 **Neural field model**

140 Upon the presentation of a stimulus, recurrent interactions between units yield complex model  
141 dynamics. In particular, population responses at any location are modulated first by their immediate  
142 (near and far) eCRFs, then by responses across the visual field as transient activity propagates  
143 through the network, until all unit responses settle into a steady-state. Such short and long-range  
144 interactions are modeled using coupled differential equations and the steady-state solution of the  
145 resulting neural field model is computed using numerical integration methods (see Materials and  
146 Methods).

147 Next, we describe experiments conducted *in silico* to compare model responses to published  
148 psychophysics data. Psychophysics studies typically record perceptual judgments related to a  
149 center stimulus under varying surround conditions. To approximate these judgments, we use an  
150 ideal neural observer which maps center population responses to a sensory value. Note that in most  
151 cases, several columns may be located within the center stimulus; while any of these columns  
152 would be suitable for readout by the ideal observer, we selected the center-most column for  
153 simplicity (unless specified otherwise). Surround modulation thus gets translated into measurable  
154 perceptual changes in the center that can then be compared to human behavioral data.

155 We have organized these experiments into three broad categories, each reflecting a key computational

156 mechanism and highlighting commonalities across visual modalities. As we will show, these  
157 experiments allow a clear picture to emerge: The diversity of observed contextual phenomena may  
158 result from a balance between two opposing “forces” that arise from complementary excitatory-inhibitory  
159 eCRF mechanisms. Figure 2 shows examples of CRF and eCRF population responses recorded  
160 from the model together with representative transformations they undergo as a result of these two  
161 forces (see Discussion for more details).

162 All model parameters (Table S1) governing the dynamics and relative strengths of the interactions  
163 between the CRF and the eCRF sub-regions were initially adjusted for the model to reproduce  
164 a host of V1 neurophysiology data (see Supplementary Experiments; Figures S1-S5) including a  
165 comparison with data from Busse et al. (2009) and Trott and Born (2015). They were held fixed for  
166 all subsequent comparisons with psychophysics data. After scaling the stimulus, the only model  
167 parameter that was optimized for individual modalities was the tuning bandwidth of individual  
168 model units.

### 169 **Competitive activation of the near vs. far surrounds**

170 Our comparison between experimental and model data starts with a set of three experiments that  
171 span the orientation, motion and color domains. All experiments involve simple center-surround  
172 stimuli, in which the surround stimulus is expected to jointly activate both the excitatory and  
173 inhibitory components of the eCRFs. Thus, these experiments should reveal a fundamental aspect  
174 of the model: the outcome of a competition between the near facilitatory and the far suppressive

175 eCRFs when they are simultaneously activated by a surround stimulus.

176 The orientation tilt occurs when the perceived orientation of a center stimulus is biased either  
177 towards (Figure 3A) or away (Figure 3B) from the orientation of a surround stimulus, also called  
178 the inducing stimulus or inducer (O'Toole and Wenderoth, 1977; Goddard et al., 2008). Figure 3C  
179 shows representative psychophysics data (digitally extracted from Figure 4 in O'Toole and Wenderoth,  
180 1977; only data averaged across subjects are available from that study). These data are characterized  
181 by two regimes: a *repulsive* regime (i.e., the perceived center orientation shifts away from the  
182 surround orientation, corresponding to positive ordinates) when the surround orientation is similar  
183 to that of the center and an *attractive* regime (i.e., the perceived center orientation shifts towards  
184 the surround orientation, corresponding to negative ordinates) when the surround orientation is  
185 different enough from that of the center.

186 The model successfully reproduces this balance between attraction and repulsion (Figure 3C; a  
187 similarly good fit was also obtained using broadband oriented textures as done in Goddard et al.,  
188 2008; data not shown). The key mechanism which enables the emergence of these two regimes  
189 is the postulated asymmetry between facilitatory and suppressive interactions originating from  
190 the near and far eCRFs, respectively. The net inhibition in the model, unlike excitation which is  
191 only dependent on pre-synaptic activity, increases monotonically with the level of post-synaptic  
192 activity of a target unit. We have confirmed this hypothesis via selective lesioning of the model  
193 key components (see Figure S14).

194 As a result, when neural population responses in the CRF and eCRF overlap significantly (as

195 when center and surround orientations are similar), inhibition predominates and center population  
196 responses get comparatively more suppressed at orientations close to that of the surround. The  
197 center of mass of center population response curves shifts away from the surround orientation,  
198 biasing the neural decoding accordingly (Figure 3A). The surround thus acts as a repellent in  
199 this regime. In contrast, when neural population responses in the CRF and eCRF are far more  
200 offset (as when the surround orientation is near orthogonal to that of the center), excitation from  
201 the near eCRF predominates, and increases the activity of center units selective for the surround  
202 orientation. This results in a force that pushes the center population response towards the surround  
203 orientation. This, in turn, biases the decoding of the center orientation in the direction of the  
204 surround orientation (Figure 3B). The surround thus acts as an attractor in this regime.

205 Beyond the orientation domain, tilt effects have also been reported for the perception of motion  
206 direction. Figure 3D shows representative psychophysics data (digitally extracted from Figure  
207 3 “periphery” condition in Kim and Wilson, 1997). Unlike in the orientation domain, however,  
208 perceptual shifts are always repulsive (the perceived motion direction of the center grating tilts  
209 away from that of the surround grating; both gratings have the same contrast and speed). This  
210 phenomenon can also be induced using coherently moving random dots (Marshak and Sekuler,  
211 1979); the effect seems to peak for similar center-surround differences in motion direction (between  
212  $40^\circ$  and  $60^\circ$ ) for either kind of stimuli.

213 We found both a qualitatively and quantitatively good fit between the model and psychophysics  
214 data as shown in Figure 3D. In the model, the disappearance of the attractive regime is accounted  
215 for by a broadening of the tuning curves (compared to orientation; see Supplementary Materials

216 and Methods). Interestingly, this seems consistent with neurophysiology data from the primary  
217 visual cortex (Ringach et al., 2002; Albright et al., 1984), which suggest that tuning for motion  
218 direction tends to be broader than for orientation.

219 In our next experiment, we show that the model is also able to account for tilt effects in the hue  
220 domain, more widely known as color induction. The model reproduces the known shifts in human  
221 judgment obtained when a center hue is surrounded by an isoluminant background of a different  
222 hue (digitally extracted from Figure 2 in Klauke and Wachtler, 2015; averaged across multiple  
223 combinations of center-surround hues sampled uniformly and independently as done in the original  
224 experiment).

225 As with motion induction, only a repulsive tilt effect is observed with hue. The model's ability  
226 to account for these data is evident from Figure 4B, which confirms the hypothesis by Klauke  
227 and Wachtler (2015) that color induction is in fact just another tilt effect (i.e., a "hue tilt effect").  
228 Furthermore, the same mechanisms that are responsible for the tilt effect in the orientation and  
229 motion domains, namely the balance between facilitatory and suppressive forces originating from  
230 the eCRF, are also at play in color induction (Figure 4A). However, opponent color coding yields  
231 populations from the center and the surround with high overlap, which explains the absence of an  
232 attractive regime for this phenomenon.

233 **Exclusive activation of the near vs. far surrounds**

234 Previous experiments involved stimuli that reflected the outcome of a competition between the near  
235 and far eCRF, which were activated jointly. Here, instead, we consider experiments that are based  
236 on surround stimuli that activate the near or the far eCRFs separately.

237 In classical depth induction experiments (Westheimer, 1986; Westheimer and Levi, 1987), human  
238 observers are presented binocularly with a center test stimulus (e.g., a thin bar) flanked by two  
239 surround stimuli (e.g., parallel thin bars or small squares). The disparities of the flanker stimuli  
240 are adjusted so that they appear in the same depth plane, either slightly in front of or behind the  
241 center stimulus. The planar separation between the center and flanker stimuli (i.e., their distance  
242 in the fronto-parallel plane) is varied systematically. Examples where the flankers appear behind  
243 the center stimulus for a shorter and a larger separation are shown in Figure 5A-B).

244 Results from the original study (data digitally extracted from Figure 1, upper panels in Westheimer  
245 and Levi, 1987) are shown in Figure 5C. When the flankers are close enough to the center stimulus,  
246 they seem to attract it in depth (corresponding to a negative shift in perceived disparity for very  
247 small flankers/center separations). That is, the center stimulus appears closer to (further away  
248 from) the observer when the flankers are in front of (behind) the center stimulus. Instead, when the  
249 flankers are moved far enough laterally, they start to repel the center stimulus (corresponding to a  
250 positive shift in perceived disparity for larger flankers/center separations).

251 The observed shifts in depth found in the model (Figure 5C) matches qualitatively with human  
252 psychophysics data: Flanker stimuli located close enough to the test stimulus activate solely the

253 near eCRF, resulting in a purely facilitatory net eCRF influence. As with the aforementioned tilt  
254 effects, net facilitatory eCRF contributions yield attraction of the center towards the surround.  
255 Conversely, flankers that are far enough from the test stimulus activate solely the far eCRF. This  
256 results in a net suppressive eCRF influence, which translates into repulsion of the center away from  
257 the surround.

258 A classical stimulus used in motion direction induction (Murakami and Shimojo, 1993; 1996) is a  
259 center-surround stimulus consisting of randomly moving dots with some coherence in the surround  
260 but no coherence in the center. The presentation of coherently moving dots in the surround elicits  
261 the illusory perception of coherent motion in the center – either in the same or in the opposite  
262 direction to that of the surround (depending on the experimental condition). In (Murakami and  
263 Shimojo, 1993; 1996), the diameter of the center and surround was fixed to  $w$  and  $2w$ , respectively,  
264 and the parameter  $w$  was varied systematically. This allowed the dimension of the overall stimulus  
265 to vary while the relative size of the center and surround regions were maintained.

266 Figure 6C shows psychophysical data (digitally extracted from Figure 6 and 7 in Murakami and  
267 Shimojo, 1996). For small stimulus sizes, the induced center movement is in the same direction  
268 as that of the surround. This corresponds to an attractive regime measured as a negative shift  
269 in the point of subjective equality (PSE). For larger sizes, the center induced movement reverses  
270 direction. This corresponds to a repulsive regime measured as a positive shift in PSE.

271 We found these results to be consistent with the model (Figure 6C). With a small enough stimulus,  
272 the coherently-moving surround dots activate the model near eCRF exclusively (Figure 6A). This

273 leads to a perceptual shift in the direction of the surround, consistent with the analogous case  
274 discussed in depth induction. At the population level, facilitation from the near eCRF tends to  
275 cause a sharpening in the population response around the surround stimulus value in an otherwise  
276 flat population response (as all motion directions are present in the center stimulus). As a result,  
277 the center stimulus looks “more like” the surround stimulus. As the stimulus size increases,  
278 the coherently-moving surround dots start to activate an increasingly large proportion of the far  
279 surround (Figure 6B), which yields the opposite (repulsive) effect. At the population level, suppression  
280 from the far eCRF causes a small notch around the surround stimulus value and the center stimulus  
281 appears to look “less like” the surround.

## 282 **Cooperative activation of the near and far surrounds**

283 Thus far, we have seen that a variety of contextual phenomena can be explained as resulting from  
284 a balance between two opposing forces: an attractive force derived from facilitatory mechanisms  
285 originating from the near eCRF vs. a repulsive force derived from suppressive mechanisms originating  
286 from the far eCRF. This competition can be tipped from attraction to repulsion by increasing  
287 the relative contribution of suppressive mechanisms originating from the far eCRF (relative to  
288 facilitatory mechanisms from the near eCRF) either by increasing the spatial extent of the stimulus  
289 (so as to activate an increasingly large proportion of the far eCRF) or by increasing the similarity  
290 between the center and surround stimulus (so as to increase the overlap between center and near  
291 surround population responses). However, we reasoned that if a surround stimulus takes on distinct  
292 and appropriate values in the near and far eCRFs (which we deem the near and far values),

293 attraction towards the near value could go in the same direction as repulsion from the far value.  
294 Thus, the joint activation of the two eCRF sub-regions would cooperate rather than compete,  
295 resulting in an even larger perceptual shift compared to what would be achieved by presenting  
296 either the near or the far stimulus values alone.

297 The implication for color perception would be that assimilation, the attraction of the perceived  
298 center hue towards a neighboring inducing hue (i.e., the near hue), could be amplified by adding an  
299 appropriate outer hue (i.e., the far hue). This idea seems consistent with an “enhanced color shift”  
300 illusion discovered by Monnier and Shevell (2003), for which we provide a novel explanation.  
301 In the classical color assimilation illusion, a colored test ring (e.g., orange) is presented within  
302 a narrow uniform surround (e.g., purple or lime), which then attracts the test ring towards its  
303 own hue. This effect was found to be greatly amplified when patterned rings (e.g., alternating,  
304 thin rings of purple and lime) at an appropriate spatial frequency and phase were used in place  
305 of the uniform colored surround (Figure 1C). Such enhancement has also been documented with  
306 achromatic stimuli (Anstis, 2006) and in brightness perception (White, 1979; Anstis, 2006).

307 Our model provides a simple explanation: As we have established, attraction (i.e., assimilation)  
308 towards say purple is caused by the activation of the near surround by a purple stimulus, with  
309 respect to a center region coinciding with the test ring. For the appropriate spatial frequency  
310 (Figure 7A), the additional lime-colored stimulus activates the far surround, leading to repulsion  
311 (i.e., contrast) *away* from lime, thus amplifying the perceptual shift towards purple as purple and  
312 lime are roughly perceptual opposites. By reversing the phase of the color grating (Figure 7B),  
313 the colors stimulating the near and far eCRFs switch, leading to the same effect in the opposite

314 direction.

315 The original psychophysics data (digitally extracted from the “6 min test” curves of Figure 5 from  
316 Shevell and Monnier, 2005) and model data are shown in Figure 7C. The model explains the  
317 existence of an optimal spatial frequency value, which maximizes the magnitude of the illusion.  
318 The spatial frequency of the stimulus controls the strength of the illusion because it determines how  
319 cleanly each of the inducing colors (e.g., lime and purple) activate the near and far eCRFs respectively  
320 for a CRF centered on the test ring. The model also postdicts that reversing the phases of the color  
321 grating leads to an effect with the same amplitude but opposite direction.

322 Critically, our explanation only depends on the appropriate hues falling within the near and far  
323 eCRFs regions; thus, we predict that the periodicity of the inducing stimulus per se is not important,  
324 as long as both regions are correctly stimulated. We show this with our own versions of the illusion  
325 in Figures S6, suggesting that the illusion is just as strong, if not stronger, when the outer rings are  
326 replaced with a single uniform region that activates the far surround optimally (which is not the  
327 case for the original stimulus by Monnier and Shevell, 2003).

## 328 **Discussion**

329 We have described a computational neuroscience model of recurrent cortical circuits to account  
330 for classical (CRF) and extra-classical receptive field (eCRF) effects. The model was constrained  
331 by anatomical data and shown in our experiments to be consistent with V1 neurophysiology.  
332 In particular, the model unifies several electrophysiology phenomena such as (cross-orientation)  
333 normalization within the CRF (Busse et al., 2009) and modulation by the eCRF (including feature-selective

334 suppression, see Trott and Born, 2015) into a computational neuroscience model of contextual  
335 integration.

336 The model further provides computational evidence for the existence of two eCRF mechanisms  
337 with complementary contributions to the CRF (a facilitatory near vs. suppressive far eCRF). In  
338 addition, the model predicts that an asymmetry between excitation and inhibition in the eCRF is  
339 needed: In our implementation, excitation depends on pre-synaptic activity only, whereas inhibition  
340 depends on both pre- and post-synaptic activities. Another model prediction is that short-range  
341 connections within a hypercolumn are weakly tuned or untuned, whereas long-range connections  
342 across hypercolumns are tuned. We ran a systematic “lesioning” study on the model, whereby each  
343 of the hypothesized mechanisms was removed individually while all remaining parameters were  
344 optimized to fit behavioral data across all experiments (see Supplementary Experiments; Figures  
345 S7-S11.)

346 Although our analysis revealed that a model which includes all assumed mechanisms performs  
347 best, we also found that some of the assumptions could be relaxed. Most importantly, a spatial  
348 segregation between the near excitatory and far inhibitory eCRF does not appear necessary and the  
349 model was found to be robust to significant overlap between these two regions (Figure S8). More  
350 generally, the model was robust to a range of parameter values (Figure S9-S11) even when relaxing  
351 the strict one-to-one mapping for the “tuned” connections from the eCRF onto the CRF (Figure  
352 S10). At the same time, while this study has focused on explaining behavioral data for an average  
353 observer, the model’s variations associated with changes in individual parameter values may help  
354 explain inter-subject variations observed experimentally (Figure S13).

355 The model distinguishes itself from previous work in succeeding to account for an array of disparate  
356 contextual phenomena spanning experimental conditions. Previous computational models have  
357 focused on explaining one or a few eCRF phenomena with an emphasis on surround suppression  
358 phenomena (see Series et al., 2003; Angelucci and Shushruth, 2013; for reviews): Phenomenological  
359 models of center-surround processing (Sceniak et al., 2001; Cavanaugh et al., 2002) and other  
360 normative models of visual coding (Coen-Cagli et al., 2012; Zhu and Rozell, 2013) have been  
361 shown to provide a good fit to single-unit contrast and size tuning responses. Recurrent network  
362 models have provided a mechanistic account for some of these phenomena (see Angelucci and  
363 Shushruth, 2013; for review) and have even led to testable predictions for single-unit electrophysiology (e.g.,  
364 Rubin et al., 2015). But, none of these models have been systematically compared to a broad and  
365 diverse set of psychophysical experiments.

366 Furthermore, our model suggests that several contextual phenomena result from not one, but two  
367 opposing forces that yield systematic distortions on center population responses: repulsion from  
368 the far suppressive eCRF vs. attraction towards the near facilitatory eCRF (see Figure 2B for  
369 representative population response dynamics). By revealing commonalities between seemingly  
370 disparate perceptual phenomena, the model has helped us establish a novel taxonomy of visual  
371 illusions: We have found that the way in which individual stimuli activate these near and far  
372 eCRFs (competitively, exclusively or cooperatively; organized by columns in Figure 8) affects the  
373 qualitative behavior of the model.

374 **A novel taxonomy of contextual phenomena**

375 Contextual stimuli that yield competitive activation of the near vs. far eCRFs were found for a  
376 set of tilt illusions including orientation (O'Toole and Wenderoth, 1977; Goddard et al., 2008),  
377 motion (Kim and Wilson, 1997) and hue (Klauke and Wachtler, 2015; also known as color induction).  
378 In these stimuli, the surround spatially overlaps with both the near facilitatory and far suppressive  
379 eCRFs – activating them both competitively. Because of the asymmetry between excitation and  
380 inhibition in the model, repulsion from the surround stimulus prevails when the center and surround  
381 population responses overlap, i.e., when the center and surround stimuli are perceptually similar.  
382 Conversely, attraction towards the surround stimulus prevails when such overlap is minimal, such  
383 as when the center and surround stimuli are perceptually dissimilar.

384 Previous authors (Klauke and Wachtler, 2015; Goddard et al., 2008; Kim and Wilson, 1997;  
385 Clifford, 2014) have suggested that surround inhibition may be key to explaining the repulsive  
386 regime in tilt effects (see Supplementary Discussion for a more in-depth discussion). The proposed  
387 mechanisms, which include shifts in neural tuning curves (Klauke and Wachtler, 2015), varying  
388 inhibition strength depending on the relative center-surround orientation (Goddard et al., 2008), or  
389 recurrent center-surround interactions (Kim and Wilson, 1997) are all consistent with the proposed  
390 mechanistic model. In addition, the present study offers a plausible computational explanation  
391 for not only the existence of a repulsive regime but also an attractive one for certain classes of  
392 stimuli, in agreement with a host of experimental data (O'Toole and Wenderoth, 1977; Goddard  
393 et al., 2008; Kim and Wilson, 1997; Westheimer and Levi, 1987).

394 Another model postdiction is the absence of such attractive regime for contextual stimuli that yield  
395 broad-band population responses (arising because of broad neural tuning for the perceptual domain  
396 or because the stimulus is inherently ambiguous as in textures with little coherent orientation). For  
397 such stimuli, the overlap between center and surround population responses remains large even  
398 for maximally dissimilar center and surround stimuli, and the only discernible contextual effect is  
399 governed by the repulsive regime. Interestingly, the model achieves its quantitative fit for motion  
400 induction experiments via a broadening of neural tuning curves for motion direction compared to  
401 orientation, which is consistent with V1 electrophysiology data (Ringach et al., 2002; Albright  
402 et al., 1984) (see also Supplementary Discussion and Figure S12 for a more in-depth discussion).

403 Stimuli that activate exclusively the near or the far eCRF have been used in classical induction  
404 experiments in the domain of depth (Westheimer and Levi, 1987) and motion (Murakami and  
405 Shimojo, 1996). In the model, consistent with the proposal by Murakami and Shimojo (1993),  
406 re-scaling a stimulus display (or similarly, varying the relative spacing between center and surround  
407 stimuli) yields a reversal from attraction to repulsion. A surround stimulus close to the center or  
408 presented at a small scale tends to predominantly activate the facilitatory near eCRF, yielding  
409 attraction towards the surround. A surround stimulus farther from the center or presented at a  
410 larger scale tends to activate the suppressive far eCRF to a greater extent, yielding repulsion away  
411 from the surround.

412 For the last set of illusions called enhanced color shifts (Shevell and Monnier, 2005), the contextual  
413 (or surround) stimulus took on “opposite” optimal values in the near and the far eCRFs. As a  
414 result, shifts induced by either region of the eCRF tended to cooperate rather than compete with

415 one another. This resulted in a perceptual shift greater than a purely attractive effect involving only  
416 the near eCRF or a purely repulsive effect involving only the far eCRF. In addition, the spatial  
417 antagonism of the model eCRF captures the existence of an optimal spatial frequency (and phase)  
418 such that a cycle of the surround stimulus coincides maximally with the near and far eCRFs.  
419 More generally, the model confirms the consensus that assimilation predominates at higher spatial  
420 frequencies and finer scales whereas contrast emerges at lower spatial frequencies and coarser  
421 scales (Murakami and Shimojo, 1993; 1996; Monnier and Shevell, 2003; Shevell and Monnier,  
422 2005; White, 1979; 1981; Anstis, 2006).

423 Shevell and Monnier (2005) have previously modeled enhanced color shifts through an S-cone  
424 color opponent model (see Supplementary Discussion). As in our model, such center-surround  
425 spatial antagonism results in the existence of an optimal spatial frequency. By design, their model,  
426 however, predicts the existence of enhanced perceptual shifts for S-cone stimuli only. One the  
427 other hand, our model predicts that such enhanced color shifts should persist for surround stimuli  
428 that do not activate S cones. We have created such stimuli (Figure S6) for the reader to judge for  
429 themselves but careful psychophysical work using a properly calibrated monitor will be needed to  
430 test this model prediction.

431 We have found a further subdivision of the above taxonomy (rows in Figure 8) based on a more  
432 detailed characterization of the center stimulus and, in particular, whether it is ambiguous (e.g.,  
433 incoherently moving random dot or achromatic stimuli) or not (e.g., high-contrast gratings and  
434 bars, highly coherent moving random dot or saturated chromatic stimuli). Unambiguous center-surround  
435 stimuli yield peaked, narrow population responses (simulation results in Figure 2B-C) across the

436 visual field. The effect of the surround on a peaked center population response is to shift its center  
437 of mass, biasing the associated decoded value accordingly (see Supplementary Discussion for a  
438 discussion of the evidence of such shifts in neurophysiology studies). The shift is either towards  
439 (attraction) or away (repulsion) from the peak of the surround population response depending  
440 on whether the net effect of the eCRF is facilitatory (Figure 2B) or suppressive (Figure 2C).  
441 Ambiguous center stimuli yield broad-band (or even flat) center population responses (Figure 2D-E).  
442 These can be distorted by a peaked surround population in two ways: a bump centered at the  
443 surround stimulus value when tuned facilitation from the eCRF prevails (Figure 2D) or a notch at  
444 the surround value when tuned suppression does (Figure 2E).

445 Table S2 shows how the literature fits in the proposed taxonomy. Note that some table entries are  
446 missing for certain visual modalities, which suggests more contextual phenomena remain to be  
447 found (e.g., cooperative shifts in orientation, which would result in an “enhanced orientation tilt”).  
448 Overall, the present study thus provides a vivid example of how computational models may help  
449 re-interpret results as well as summarize and integrate disparate phenomena.

#### 450 **Open questions**

451 The neural tuning curves considered in this work (orientation, disparity, motion direction, color  
452 opponent) can be found in relatively low-level areas of the visual cortex, such as V1, V2 or MT.  
453 Thus, the consistency between model and behavioral data is all the more remarkable as many of  
454 the illusions studied here are likely to also involve higher-level visual processes which are known

455 to affect perception including perceptual organization and grouping (e.g., Manassi et al., 2016),  
456 attention and other top-down feedback (Gilbert and Li, 2013) including surface-based and other  
457 filling-in processes (Grossberg and Todorović, 1988). The model's ability to account for contextual  
458 interactions may be limited to the relatively simple stimuli such as the bars and gratings tested  
459 here. We expect the model to fail to account for human data for more complex contextual stimuli  
460 defined by objects or shapes (e.g., Manassi et al., 2016). At the very least, a more complete  
461 model would likely require multiple stages of processing as well as mechanisms of filling-in and  
462 contour extraction (Grossberg and Todorović, 1988). Similarly, considering tuning curves found in  
463 higher-level areas, such as tuning to hue observed in V4/PIT neurons (as opposed to color-opponent  
464 V1 neurons considered here, see Conway et al., 2007) could also improve the fit with experimental  
465 data (though hue tuning remains controversial, see Mollon, 2009; Conway, 2009).

466 More generally, the present model leaves open any role for attention. Indeed, recent work has  
467 shown that attention seems to be shifting both the CRF and eCRF independently towards the  
468 attended location (Anton-Erxleben et al., 2009). It is likely that attention (not accounted for in  
469 the present model) may have played a role in shaping the pattern of observed behavioral results. In  
470 our simulations, the CRF size was scaled to the center of the stimuli – a role that could possibly be  
471 endowed to attention (Carandini, 2012). Indeed, one of the main mechanisms in the present model –  
472 that of complementary excitatory and inhibitory surround mechanisms – is a key mechanism in one  
473 of the leading models of spatial attention (Tsotsos et al., 2001). In this model, an annular region  
474 of inhibition creates a negative attentional field surrounding the region of perceptual facilitation  
475 centered on the attended target. In addition, modeling work has also suggested that top-down

476 influences may “gate” the effective contextual interactions mediated by long-range horizontal  
477 connections (Setić and Domijan, 2008).

478 Anatomical data to constrain the patterns of recurrent connectivity (both within and across hypercolumns)  
479 in the model are scarce. The near and far eCRFs as modeled are likely to constitute, at best,  
480 coarse approximations for more complex patterns of anatomical connections. In particular, both  
481 the spatial extent and the relative strength of the near and far eCRFs relative to that of the CRF  
482 were held constant across experiments. Given that the experiments considered throughout spanned  
483 a range of visual stimuli across modalities and sizes, it is likely that these phenomena recruit neural  
484 populations in different cortical areas and visual eccentricities. It is also likely that variations in  
485 experimental factors lead to differences in how the center and surround capture attention. We thus  
486 expect improvements in the model’s quantitative fit by considering additional parameters to control  
487 the spatial extent and the relative strength of the near and far eCRFs (e.g., as done in Goddard et al.,  
488 2008).

489 We have also left open the question of whether the connectivity in the near and far eCRFs would  
490 draw on slow intra-areal lateral connections or fast intra-areal feedback connections (Angelucci  
491 et al., 2002a;b; Shushruth and Ichida, 2009). Such lack of refinement, in addition to a lack of  
492 realistic modeling of excitatory and inhibitory synapses and their relative timing (Vinck et al.,  
493 2013), negatively impacts our ability to make predictions about the precise time course of the  
494 contextual effects modeled. We also expect that a resolution on the question of feedback vs. lateral  
495 connectivity will be needed to account for some of the electrophysiology phenomena we left aside  
496 in the present study including the known contrast dependence of the eCRF size (see Angelucci and

497 Shushruth, 2013; for review) or cross-orientation enhancements (Levitt and Lund, 1997; Sillito  
498 et al., 1995). Another hypothesis of the model that has yet to be confirmed is the existence of  
499 cortical columns for all visual domains beyond orientation (see Sincich and Horton, 2005; for  
500 review). At present, the existence of cortical columns for color (Dow, 2002), motion (DeAngelis  
501 et al., 1999) and binocular disparity (DeAngelis and Newsome, 1999) is only partially supported  
502 by neurophysiology evidence.

503 We have assumed for simplicity that the near eCRF is circular (i.e., isotropic with respect to the  
504 topography of the visual field). There is, however, evidence for anisotropies in the pattern of  
505 horizontal connections between cortical columns (as orientation-tuned cells tend to be more often  
506 connected when they share the same selectivity and their CRFs are aligned along an axis parallel  
507 to their preferred orientation, see Bosking et al., 1997). There is also more direct evidence for  
508 anisotropies in the shape of the eCRF (i.e., various elongations over a wide range of orientations  
509 and widths, see Tanaka and Ohzawa, 2009). The function of these anisotropies has been attributed  
510 to the computation of higher-order features (including contrast- or texture-defined boundaries) as  
511 well as contour integration and pop-out (Stemmler et al., 1995; Hess et al., 2003; Tanaka and  
512 Ohzawa, 2009). Future work should test whether these phenomena can be accounted for with a  
513 model extension that incorporates such eCRF anisotropies.

514 More generally our study did not address the role of the perceptual biases and the altered discriminability  
515 that arise because of surround mechanisms. It has been suggested that surround mechanisms  
516 could constitute one of the primary mechanisms for predictive coding and Bayesian inference  
517 type of computations (see Schwartz et al., 2007; for review). We speculate that the computational

518 mechanisms revealed by the contextual illusions studied here play a key role in shaping invariant  
519 population codes for object constancy. We have obtained preliminary results suggesting that tuned  
520 suppression from the far eCRF may improve the accurate decoding of surface reflectances across  
521 changes in illumination (i.e., color constancy; see Mély & Serre, abstract presented at the 2015  
522 Vision Science Society meeting), by helping to discount undesirable variations in center population  
523 responses caused by changes in the light source. (This is reminiscent of a color constancy algorithm  
524 by Land and McCann (1971) known as the Retinex.) This raises the intriguing possibility that at  
525 least some of the mechanisms unraveled here may support other forms of perceptual constancy  
526 beyond color. Further work will be needed to quantify how object transformations such as changes  
527 in illumination or depth affect neural population responses tuned to orientation or binocular disparity  
528 and what computational mechanisms are needed to help discount these nuisances. Nonetheless,  
529 the ability of the model to account for the variety and complexity of contextual illusions provides  
530 computational evidence for a novel canonical cortical circuit shared across visual modalities.

## 531 **Materials and Methods**

532 Additional methods may be found in Supplemental Materials and Methods.

### 533 **Model connectivity**

534 A column centered at location  $(x, y)$  contains a complete set of  $N$  units with CRFs centered  
535 at  $(x, y)$  and tuning values covering the full range  $\theta_{k=1\dots N}$  (e.g., orientation tuning curves are  
536 regularly centered at values  $\theta_k \in [0, 180^\circ]$ ). Tuning curves are idealized – either bell-shaped for

537 disparity (Cumming and Parker, 1997), motion direction (Albright et al., 1984) and orientation (Ringach  
 538 et al., 1997) or monotonic for color opponency (Johnson et al., 2001).

539 Each unit  $(x, y, k)$  receives excitation  $Q_k^{xy}$ , assumed to be weakly tuned and originating from within  
 540 the same hypercolumn:

$$Q_k^{xy} = \sum_{j=1\dots N} w_{jk} X_j^{xy} \text{ s.t. } w_{jk} \sim \mathcal{N}(\theta_k, \zeta), \quad (1)$$

541 where  $w_{jk}$  corresponds to excitatory weights between units  $k$  and  $j$  (with selectivity  $\theta_k$  and  $\theta_j$ ,  
 542 respectively) and  $X_j^{xy}$  to input activity at location  $(x, y, j)$ . We assume these weights to be normally  
 543 distributed, centered at a target unit tuning preference  $\theta_k$  with standard deviation  $\zeta$ . Some tuning  
 544 (albeit weak) is necessary in order to prevent intra-columnar excitation from flattening the population  
 545 responses to well-defined stimuli. In the color domain, we consider color-opponent model units  
 546 with monotonic tuning curves. Instead of drawing weights from a normal distribution, which only  
 547 makes sense for bell-shaped tuning curves, we set  $w_{kk} = (\zeta\sqrt{2\pi})^{-1}$  and  $w_{jk} = \text{const.}$  (when  $j \neq k$ ;  
 548 under the constraint that the weights sum up to  $r1$ ).

549 Each unit  $(x, y, k)$  also receives some inhibition  $U^{xy}$ , assumed to be untuned and originating from  
 550 within the same hypercolumn:

$$U^{xy} = \frac{1}{N} \sum_{j=1\dots N} Y_j^{xy}, \quad (2)$$

551 where  $Y_j^{xy}$  is the output activity of unit  $j$  at location  $(x, y)$ . Unlike the excitation which is linear,

552 inhibition is non-linear, of the shunting kind (Grossberg and Todorović, 1988), and acts on the  
553 output of the pre-synaptic units (Equation 5). Combined with broad tuning, this allows populations  
554 of cells with coinciding CRFs to be significantly driven by a common input. Such mechanism was  
555 used in the model by Shushruth et al. (2012) and found experimentally to be critical to reproduce  
556 nonlinear neural effects such as stimulus-matched surround suppression (Trott and Born, 2015).  
557 Because the local inhibition is untuned, its strength is independent of a unit selectivity  $\theta_k$ , and we  
558 drop the subscript  $k$  for simplicity.

559 Furthermore, unit  $(x, y, k)$  also receives tuned excitation  $P_k^{xy}$  from other units with the same selectivity  
560  $\theta_k$  that are located within its near eCRF  $\mathbb{N}^{xy}$ , defined relatively to position  $(x, y)$ :

$$P_k^{xy} = \frac{1}{|\mathbb{N}^{xy}|} \sum_{u,v \in \mathbb{N}^{xy}} X_k^{u,v} \quad (3)$$

561 Similarly, unit  $(x, y, k)$  also receives tuned inhibition  $T_k^{xy}$  from other units with the same selectivity  
562  $\theta_k$  that are located within its far eCRF  $\mathbb{F}^{xy}$ , defined relatively to position  $(x, y)$ :

$$T_k^{xy} = \frac{1}{|\mathbb{F}^{xy}|} \sum_{u,v \in \mathbb{F}^{xy}} Y_k^{u,v} \quad (4)$$

563 As in Equation 2, inhibition is non-linear and acts on the output  $Y_k^{u,v}$  of unit  $k$  at location  $(u, v)$ .

564 **Neural field model**

565 Neural field dynamics obey the following equations:

$$\begin{aligned} \eta \partial_t X_k^{xy} + \varepsilon^2 X_k^{xy} &= [\xi L_k^{xy} - (\alpha X_k^{xy} + \mu) U^{xy} - (\beta X_k^{xy} + \nu) T_k^{xy}]_+ \\ \tau \partial_t Y_k^{xy} + \sigma^2 Y_k^{xy} &= [\gamma P_k^{xy} + \delta Q_k^{xy}]_+ \end{aligned} \quad (5)$$

566 where the feed-forward input  $L_k^{xy}$  drives every unit  $(x, y, k)$  across the visual field; each is represented  
 567 by its recurrent input  $X_k^{xy}$  and output  $Y_k^{xy}$ . The parameters  $\alpha, \beta, \delta, \gamma, \mu$  and  $\xi$  can be interpreted as  
 568 synaptic weights (see Table S1 for values used) which control the amount of intra- and inter-columnar  
 569 excitation and inhibition (Equations 1–4). The steady-state solution is computed using numerical  
 570 integration (with convergence typically taking  $\sim 50$  iterations). Population responses at the steady-state  
 571  $\bar{Y}_k^{xy}$  are a very nonlinear function of the model input  $L_k^{xy}$ .

572 For each unit, the steady-state input and output are given by  $\bar{X}_k^{x,y}$  and  $\bar{Y}_k^{x,y}$ , resp. Due to the  
 573 rectifying non-linearity in the dynamics (Equation 5), at steady-state,  $\bar{X}_k^{x,y}$  and  $\bar{Y}_k^{x,y}$  can either be  
 574 equal to zero, or to the values below:

$$\begin{aligned} \bar{X}_k^{x,y} &= \frac{\xi L_k^{x,y} - \mu \bar{U}^{x,y} - \nu \bar{T}_k^{x,y}}{\varepsilon^2 + \alpha \bar{U}^{x,y} + \beta \bar{T}_k^{x,y}} \\ \bar{Y}_k^{x,y} &= \frac{\gamma \bar{P}_k^{x,y} + \delta \bar{Q}_k^{x,y}}{\sigma^2}, \text{ with:} \end{aligned} \quad (6)$$

$$\begin{aligned}
\bar{U}^{x,y} &= \frac{1}{N} \sum_{j=1\dots N} \bar{Y}_j^{x,y} \\
\bar{T}_k^{x,y} &= \frac{1}{|\mathbb{F}^{x,y}|} \sum_{u,v \in \mathbb{F}^{x,y}} \bar{Y}_k^{u,v} \\
\bar{P}_k^{x,y} &= \frac{1}{|\mathbb{N}^{x,y}|} \sum_{u,v \in \mathbb{N}^{x,y}} \bar{X}_k^{u,v} \\
\bar{Q}_k^{x,y} &= \sum_{j=1\dots N} w_{j,k} \bar{X}_j^{x,y}
\end{aligned} \tag{7}$$

## 575 **Tuning curves**

576 The model constitutes an example of tuning curve population model (Schwartz et al., 2007; Rust  
577 et al., 2006). We considered two kinds of tuning curves: bell-shaped (orientation, motion direction,  
578 binocular disparity) and monotonic, non-saturating tuning curves (color). All tuning curves were  
579 normalized, i.e., the maximum unit activity was set to be equal to 1. For non-angular variables  
580 (e.g., disparity), bell-shaped tuning curves were parametrized as Gaussian functions:

$$f(\theta \mid \theta_k, \sigma) = \exp\left(-\frac{(\theta - \theta_k)^2}{2\sigma^2}\right), \tag{8}$$

581 with preferred stimulus value  $\theta_k$  and tuning bandwidth  $\sigma$ . When the variable was circular (e.g.,  
582 orientation, motion direction), we modeled the tuning curve as a von Mises function instead:

$$f(\theta \mid \theta_k, \sigma) = \exp\left(\frac{\cos\left(\left(\theta - \theta_k\right)\frac{2\pi}{I}\right) - 1}{2\sigma^2}\right) \tag{9}$$

583 where  $I$  indicates the length of domain of the tuning curve (e.g.,  $\pi$  for orientation vs.  $2\pi$  for  
584 direction). We generally sampled on the order of 30 tuning curve centers regularly spaced in  
585 the domain of the considered visual modality. We found that the number of tuning curve centers  
586 considered did not impact our results as long as it was large enough.

587 Monotonic, non-saturating tuning curves for color were derived by converting stimuli to idealized  
588 cone responses first, which were then mapped to opponent color channels similarly to Zhang et al.  
589 (2012). These included red-on/green-off ( $R^+G^-$ ), green-on/red-off ( $G^+R^-$ ), blue-on/yellow-off  
590 ( $B^+Y^-$ ), and yellow-on/blue-off ( $Y^+B^-$ ), alongside with a pair of luminance-sensitive channels,  
591 selective for lighter ( $Wh^+Bl^-$ ) and darker ( $Bl^+Wh^-$ ) stimuli.

## 592 **Model parameters**

593 All circuit parameters were held constant in all comparisons with psychophysics data. They  
594 were determined a priori in order to reproduce key neurophysiology data (see Supplementary  
595 Experiments) and were held constant for all visual modalities except color because of a qualitative  
596 difference in tuning curve (see Equation 1). In all subsequent experiments, only two variables were  
597 allowed to vary: the stimulus scale and the tuning bandwidth for model units.

598 The stimuli used in psychophysics studies varied greatly – recruiting neural populations subtending  
599 a wide range of CRF (and eCRF) sizes and eccentricities, possibly spanning different visual areas.  
600 Rather than adjusting the size of the model CRFs and eCRFs for individual experiments, which  
601 would have required structural changes to the model, we instead varied the stimulus scale. Because

602 the connectivity between model hypercolumns was held fixed, this is somewhat akin to varying the  
603 magnification factor in the model. Critically, this yielded broad estimates for CRF (and eCRF)  
604 sizes within a biologically realistic range (from a fraction of a degree of visual angle to a couple  
605 of degrees). The width of the idealized tuning curves which is common to all model units was  
606 optimized separately for each experiment (see Supplementary Materials and Methods for details).  
607 We have confirmed that our key model predictions were robust over a range of these parameter  
608 values.

### 609 **Ideal neural observer model**

610 We used an ideal neural observer model to map model population responses to decoded sensory  
611 variables, which can then be compared to behavioral judgments collected experimentally. We used  
612 a population vector model (Georgopoulos et al., 1986), in which each unit votes for its preferred  
613 sensory value in proportion to its activity (normalized by the summed activities of all units within  
614 the same column). This model is not appropriate for color because of the tuning along opponent  
615 color pairs rather than a hue angle. Instead, we used cross-validated ridge regression to decode the  
616 sine and cosine of hue.

### 617 **Acknowledgments**

618 This work was supported by DARPA young faculty award [grant number N66001-14-1-4037] and  
619 NSF early career award [grant number IIS-1252951] to TS.

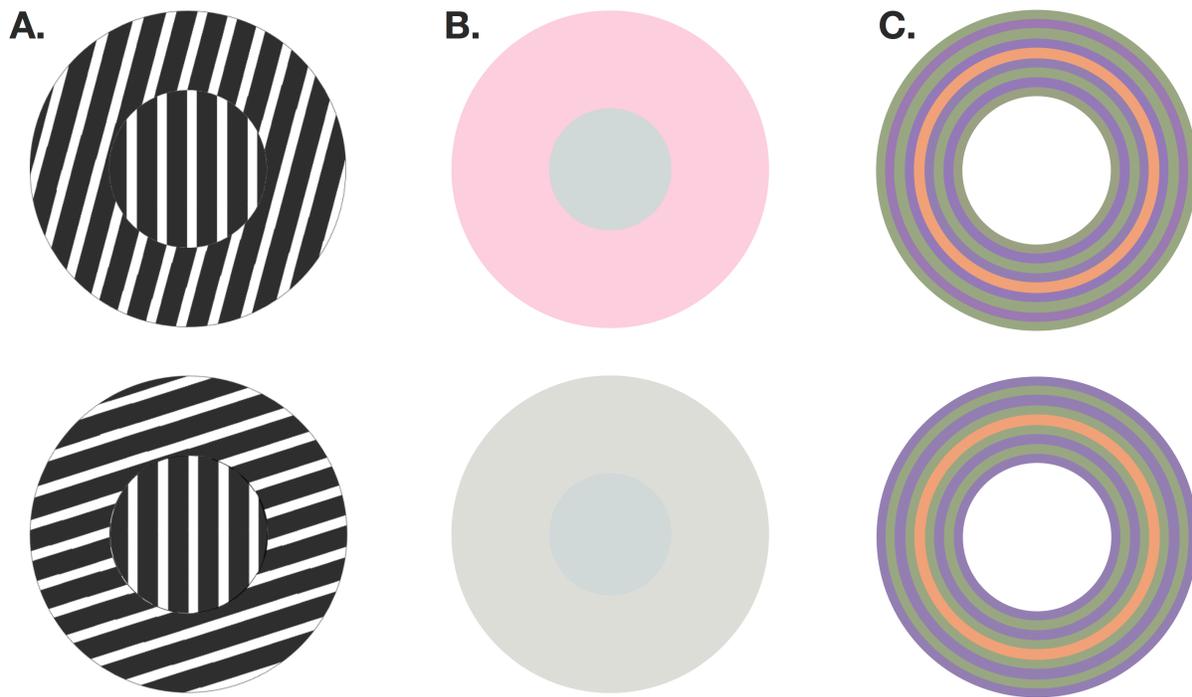


Figure 1. *Representative contextual phenomena explained by the model.* **A. Orientation tilt:** The perceived center (or test) orientation appears tilted from its true physical orientation, away from the surround (or contextual) orientation when center and surround stimuli are similar (top) and towards the surround orientation when they are dissimilar (bottom). **B. Color induction:** A central gray stimulus appears greener when embedded in a pink surround (top) compared to a neutral gray surround (bottom). **C. Enhanced color shifts:** The test stimulus is a central, orange ring, embedded in a surround stimulus composed of alternating purple and lime rings. The test ring looks vividly more pink when the adjacent color is purple, followed by lime (top), and looks more yellow when lime is the adjacent color, followed by purple (bottom).

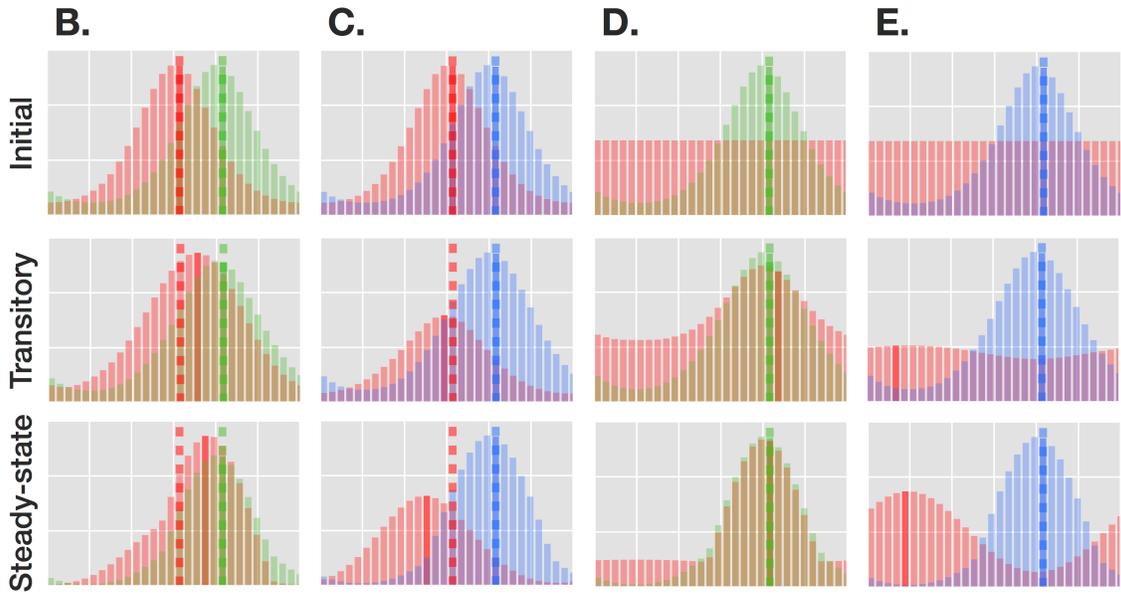
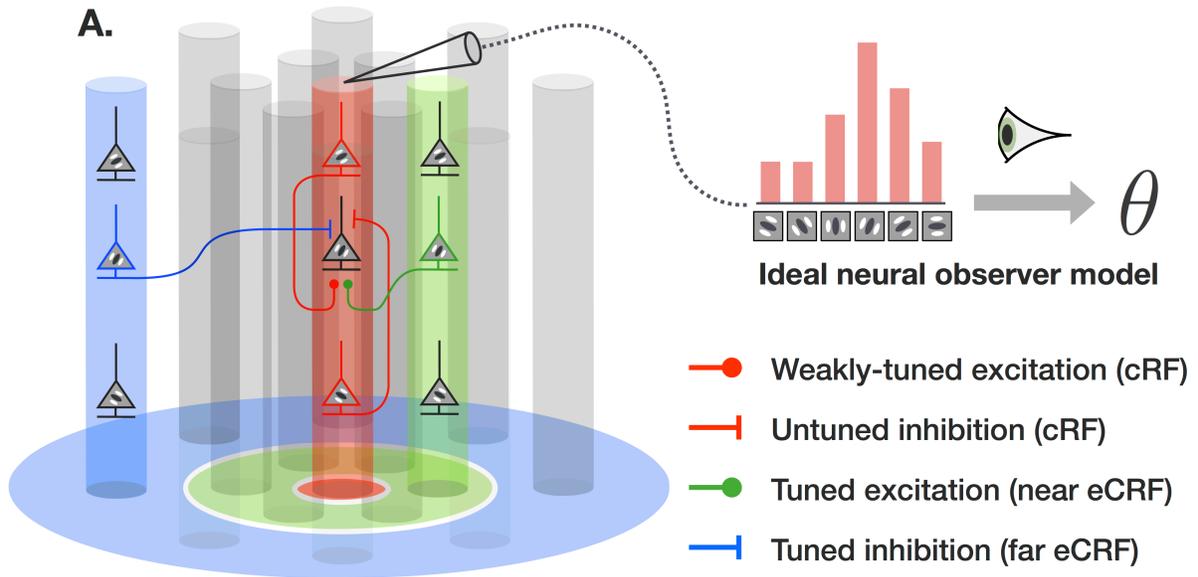


Figure 2. *Recurrent network model of center-surround interactions.* **A. Connectivity:** The model implements excitatory and inhibitory connections both short-range (within hypercolumns) and long-range (between hypercolumns). The regions shown in red, green and blue correspond to the CRF, near eCRF and far eCRF, respectively (defined for the reference column in red). Model inhibitory connections are such that the net inhibition onto a target unit is a function of not just the pre-synaptic activity, but also the post-synaptic activity (see text for details). Color conventions for the CRF and the near and far eCRFs are used consistently throughout the paper.

**B–E. Representative model dynamics:** Example population responses (32 direction-tuned model units) following the presentation of a contextual stimulus corresponding to the initial, transitory and steady state (rows). Population responses correspond to locations in the CRF, near eCRF and far eCRF. Highlighted bars represent directions decoded from the corresponding populations (undefined for flat responses); dashed lines represent initial decoded values at stimulus onset. Each column corresponds to a representative transformation undergone by the center population under the proposed taxonomy of contextual phenomena derived from the model: **B.** attractive shift, **C.** repulsive shift, **D.** bump, **E.** notch. (see also Discussion and Figure 8). Abscissas span the range  $[-180^\circ, 180^\circ]$  and ordinates are normalized independently for readability.

620

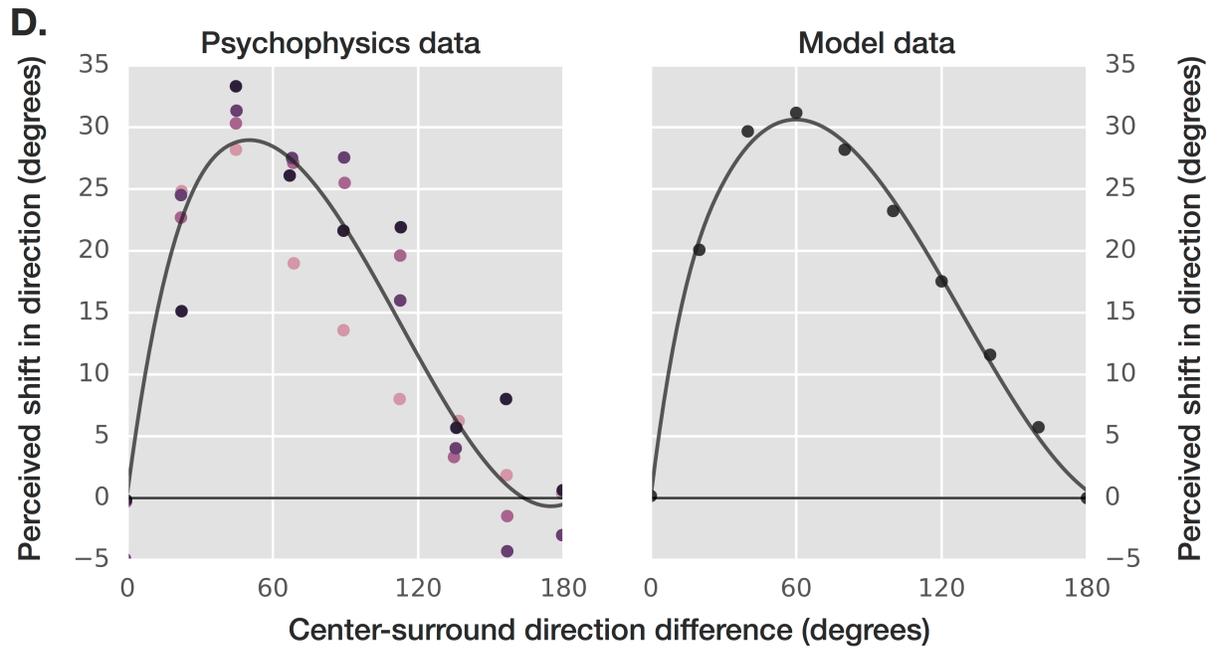
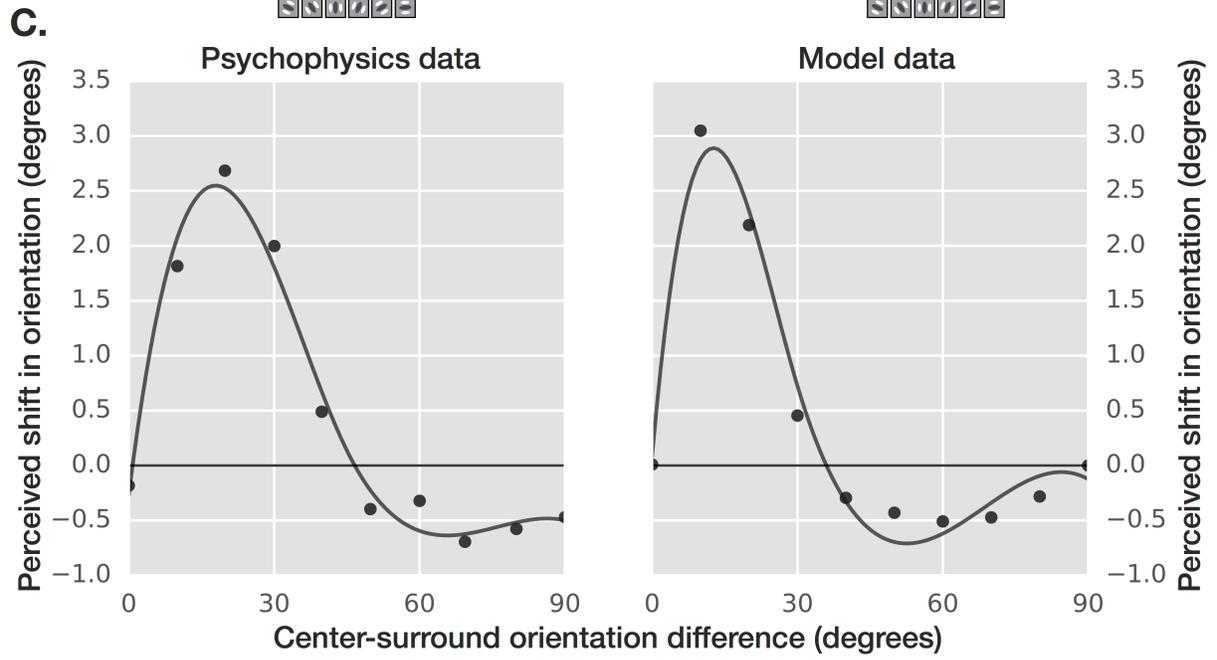
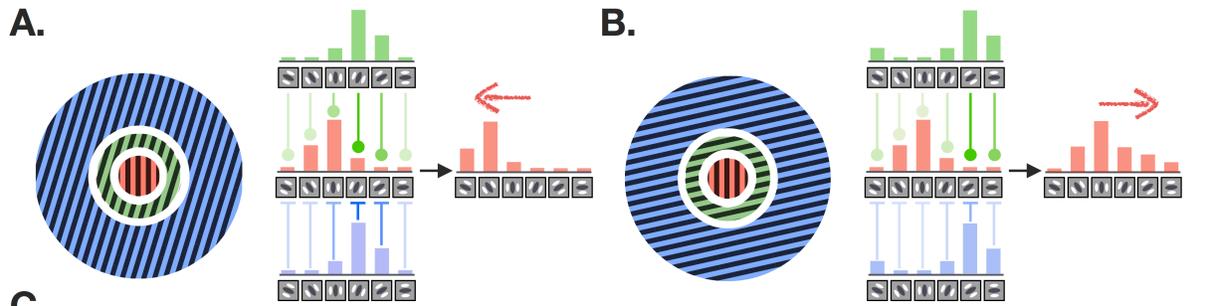


Figure 3. *Tilt effects*. Competitive activation of the near vs. far eCRFs explains the shift in tilt from one direction to the other. **A. Repulsion:** For similar center-surround orientations, tuned inhibition from the far eCRF outweighs excitation from the near eCRF, which yields a net repulsive force on the center population responses (away from that the surround orientation). **B. Attraction:** For dissimilar center-surround orientations, tuned excitation from the near eCRF prevails, which yields a net attractive force on the center population responses (towards the surround orientation). Note that gaps between the CRF and the near and far eCRFs were added for improved readability only and are not present in the actual model. **C. Orientation tilt: Psychophysics vs. model data.** Psychophysics data were digitally extracted from Figure 4 in (O’Toole and Wenderoth, 1977) and fitted with splines. The model explains the characteristic shift from perceptual repulsion (positive ordinates) to attraction (negative ordinates). **D. Motion tilt: Psychophysics vs. model data.** Psychophysics data were digitally extracted from Figure 3 (“periphery” condition) in (Kim and Wilson, 1997) and fitted with splines. Different colors correspond to different subjects. Both psychophysics and model data exhibit a similar dependency on the direction difference between center and surround, as well as a lack of an attractive regime.

621

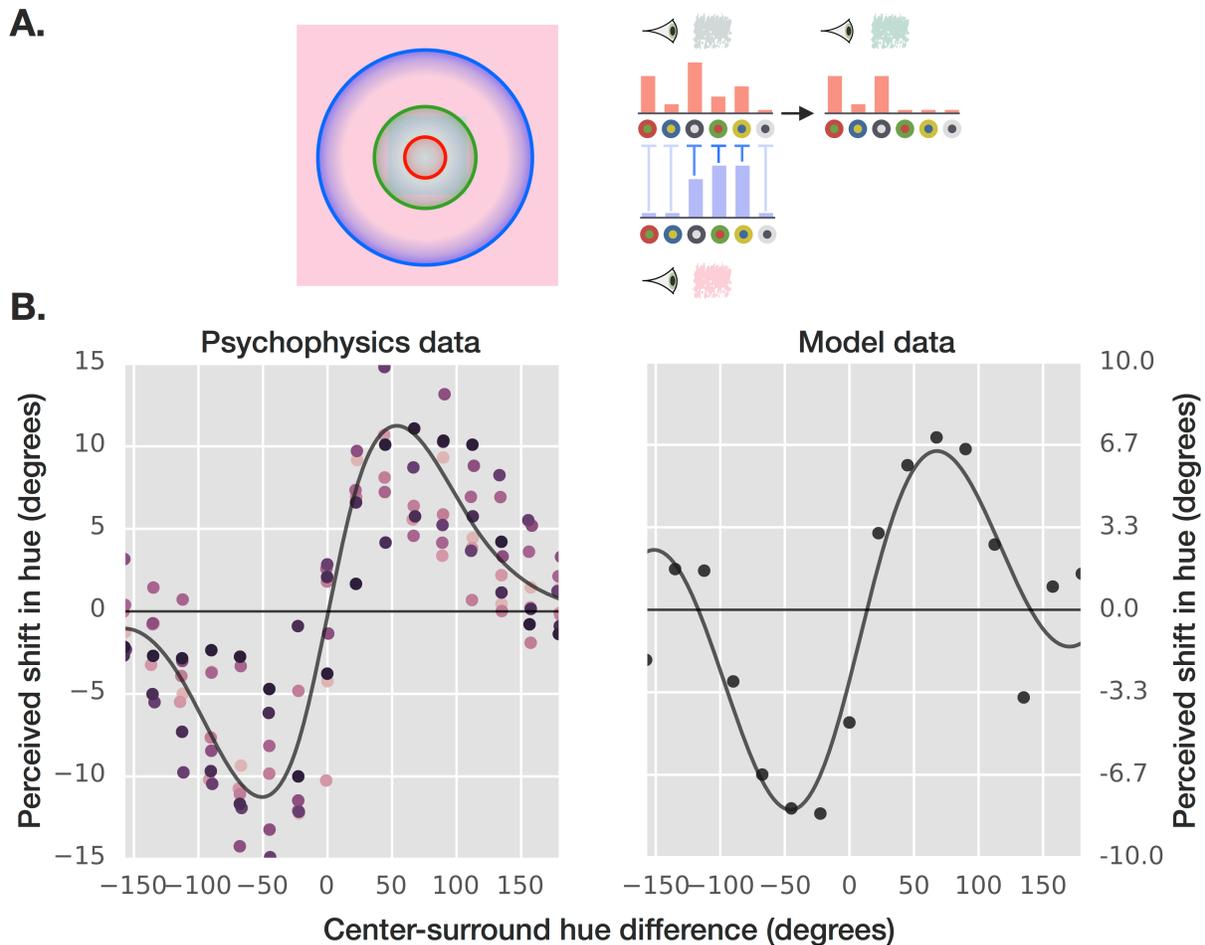
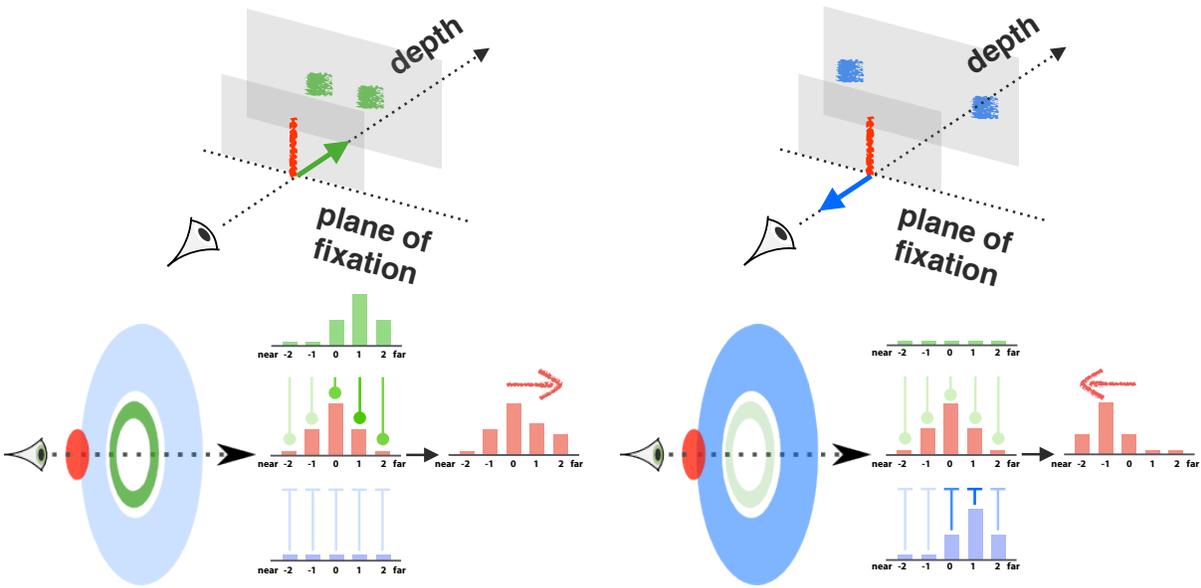


Figure 4. *Color induction (or hue tilt effect)*. This experiment generalizes the tilt effect to opponent population codes. **A. Repulsion:** As with the classical tilt effect, the key model mechanism behind perceptual repulsion is the tuned inhibition from the far eCRF. In this example, the pink surround suppresses “red” center neurons, therefore reducing the “redness” of the gray center patch yielding a shift in the perceived center hue towards green. The same explanation also applies to chromatic center stimuli. Colored patches shown next to the eyes correspond to the color decoded under the ideal observer. **B. Psychophysics vs. model data:** Psychophysics data were digitally extracted from Figure 2 in (Klaube and Wachtler, 2015) and fitted with splines (averaged across eight surround hues). Both model and behavioral data exhibit a characteristic two-lobed shape peaking around  $\pm 50^\circ$ . Please note the difference in ordinate scale between the psychophysics and model data.



C.

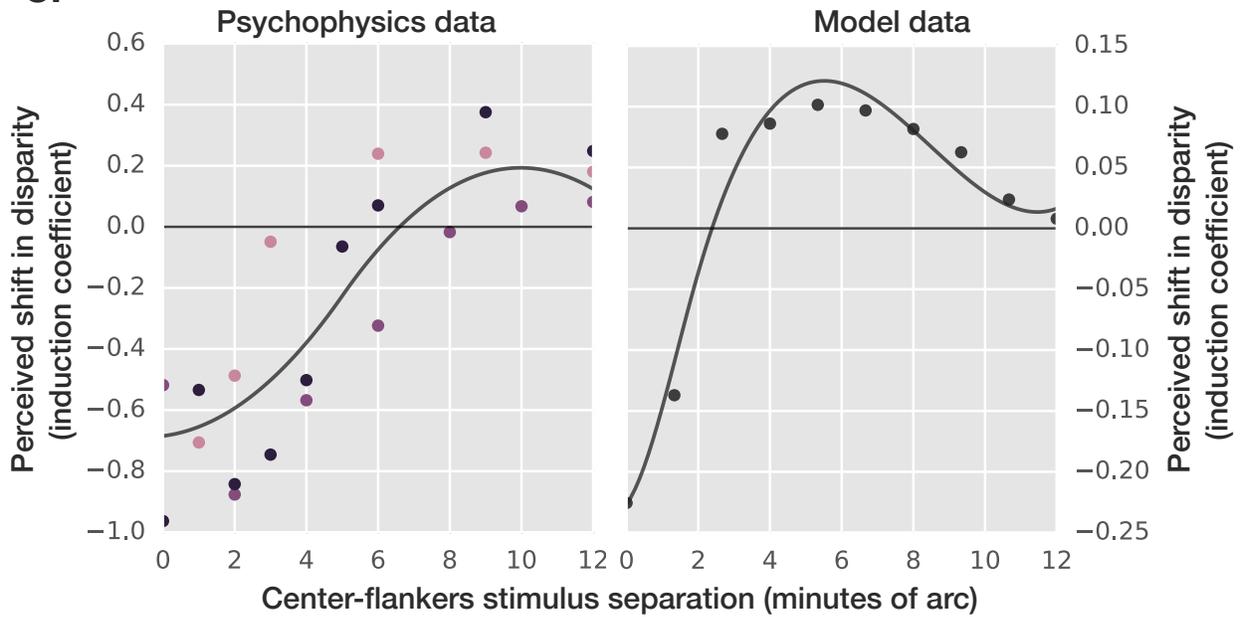


Figure 5. *Depth induction*. The exclusive activation of either the near or far eCRFs by flankers (as their separation vary) explains the existence of the shift from assimilation to contrast. The perceived depth of a binocular center stimulus (at zero disparity) is affected by binocular flankers located on either sides, presented at either crossed or uncrossed disparities. **A. Attraction:** For short separations, flankers activate the near eCRF, which yields a net attractive force on the center population responses (towards the surround disparity) corresponding to a negative perceived shift. **B. Repulsion:** For larger separations, flankers activate the far eCRF, which yields a net repulsive force on the center population responses (away from the surround disparity) corresponding to a positive perceived shift. **C. Psychophysics vs. model data:** Psychophysics data were digitally extracted from Figure 1 (upper panels) in (Westheimer and Levi, 1987) and fitted with splines. Both behavioral and model data capture the balance between stronger attraction towards the flankers at small separations, and weaker repulsion at larger separations. Note that the agreement between the model and human data is only qualitative as the perceived shifts in disparity are on different scales (the model underestimates the strength of the attractive regime in this illusion).

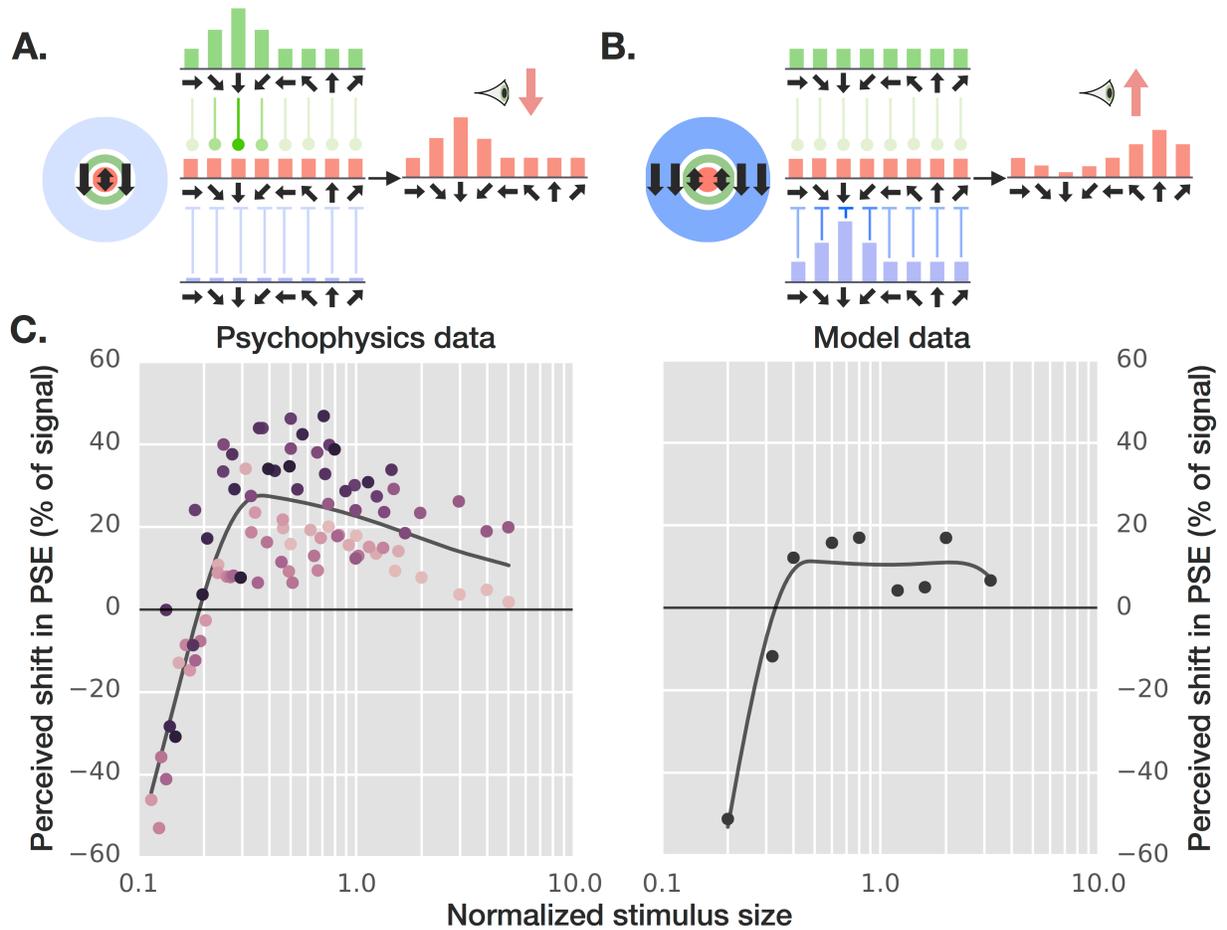


Figure 6. **Motion induction.** The increasing proportion of the far eCRF activated by larger stimuli explains the shift from assimilation to contrast. **A. Attraction:** When the overall stimulus is small enough, the coherently moving surround dots activate the near eCRF exclusively, leading to motion assimilation (i.e., the center dots' direction appears the same as that of the surround dots). **B. Repulsion:** Beyond a critical size, activation of the far inhibitory eCRF prevails, which leads to the opposite motion contrast effect (i.e., the center dots' direction appears opposite to that of the surround dots). **C. Psychophysics vs. model data:** Psychophysics data were digitally extracted from Figure 5 and 6 (Murakami and Shimojo, 1996) and fitted with splines. Both exhibit stronger attraction (negative ordinates) for smaller stimulus sizes, and weaker repulsion (positive ordinates) for larger sizes. Shifts in the point of subjective equality (PSE) were used as a proxy for shifts in perceived motion direction.

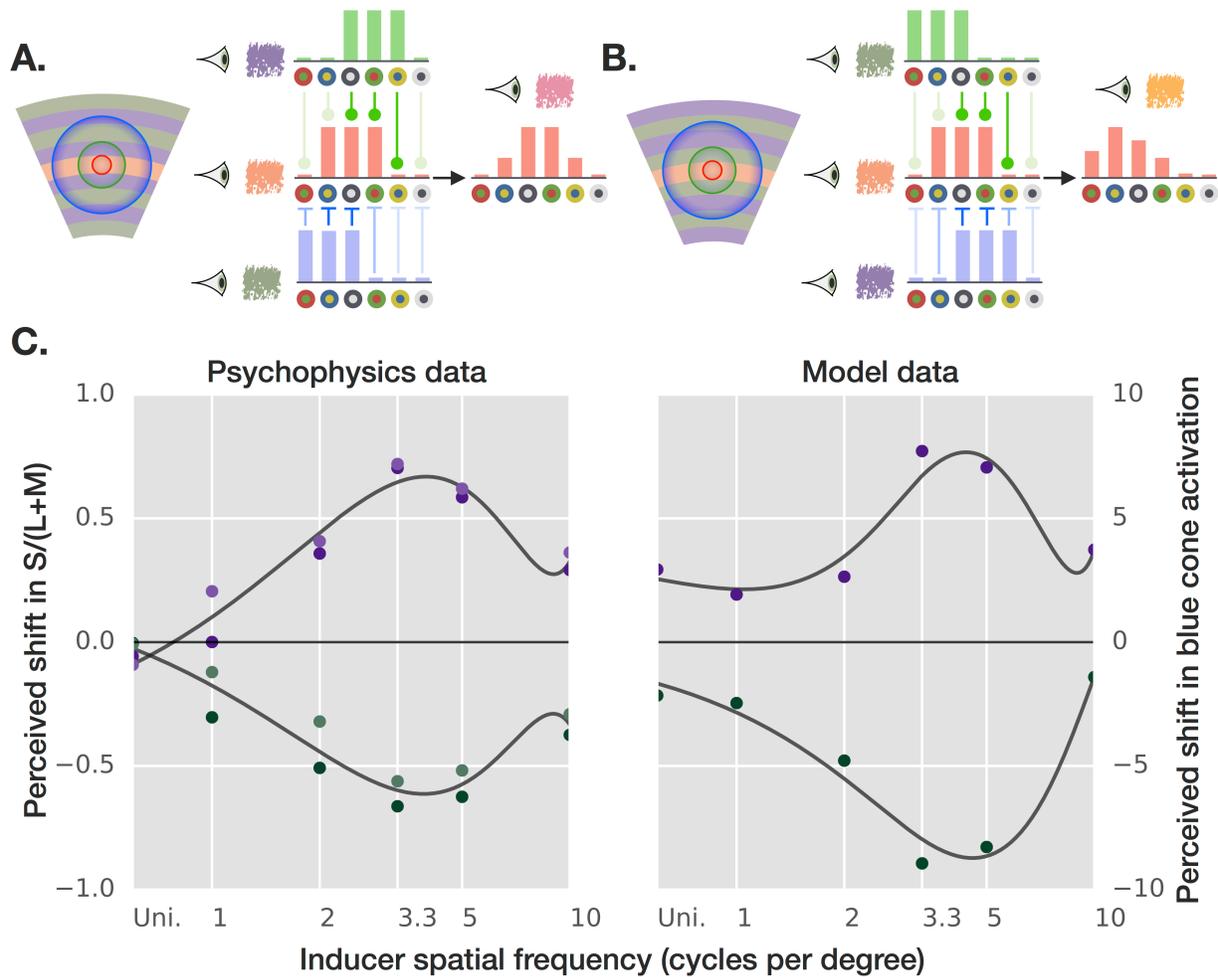


Figure 7. **Enhanced color shifts.** Cooperative activation of the near and far surrounds explains enhanced perceptual shifts. When distinct and “opposite” hues are used in a patterned surround (or inducer), the resulting shift in color perception of a test hue (here, orange) is amplified relative to a uniform surround of either hue. **A. Shift in one direction:** For the optimal spatial frequency, one surround hue (e.g., purple) overlaps optimally with the near eCRF and the other one (e.g., lime) with the far eCRF. For the right color combination (as here with purple and lime which are complementary colors), this results in cooperating perceptual forces: a shift towards purple / away from lime. The colored patches next to the eyes correspond to the color decoded under the ideal observer. **B. Shift in the other direction:** when purple and lime are switched. **C.**

**Psychophysics vs. model data:** Psychophysics data were digitally extracted from Figure 5 (6 minutes test condition) in (Shevell and Monnier, 2005) and fitted with splines. Purple/green dots correspond to condition A/B. ‘Uni.’ stands for a uniform inducer composed of a single hue. For

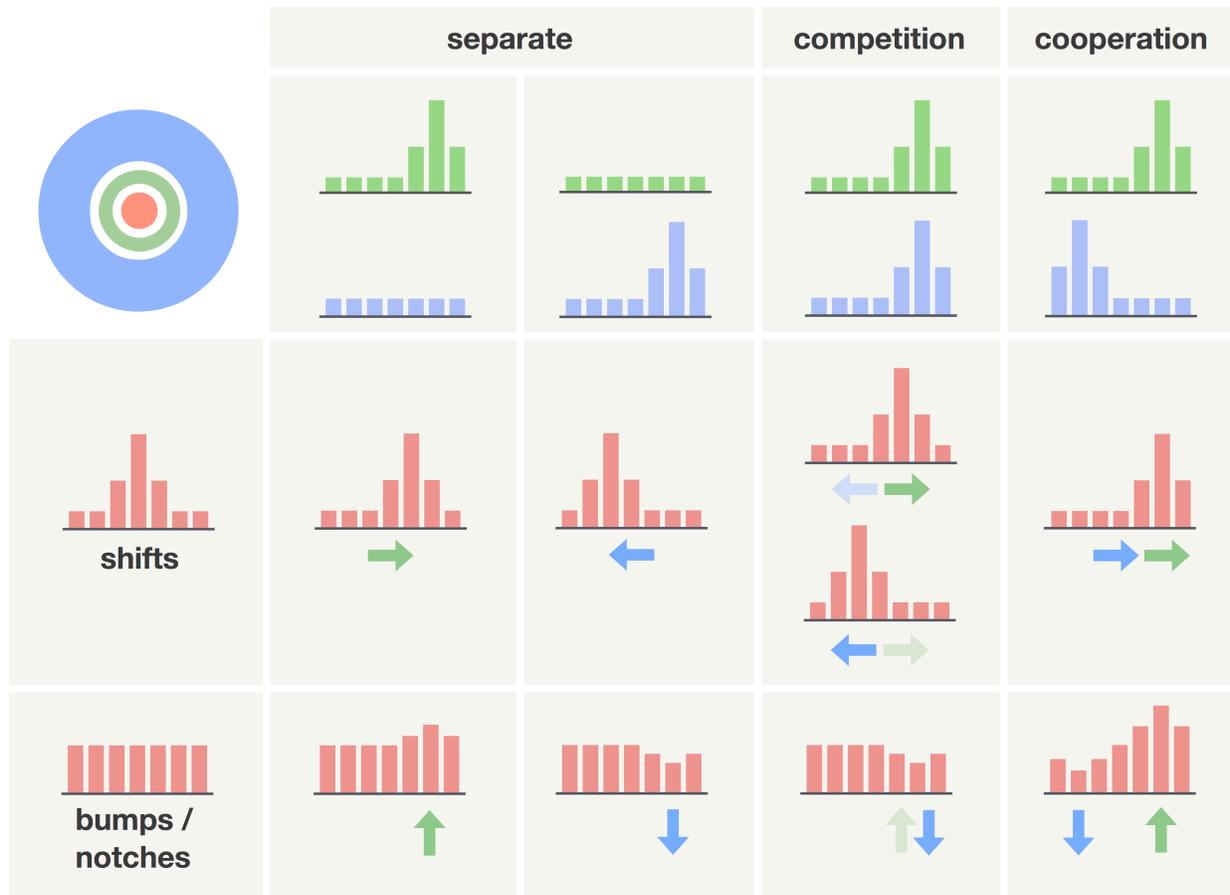


Figure 8. *A new taxonomy of contextual phenomena.* **Rows:** Contextual phenomena manifest themselves in the model either as (i) shifts with peaked center population response curves (unambiguous stimuli), or (ii) bumps/notches with broad/uniform center population response curves (ambiguous stimuli). **Columns:** Center-surround stimuli activate the near and far eCRFs in three typical ways: (i) either one separately, (ii) both competitively (i.e., near and far eCRFs each induce shifts that tend to stymie each other; green and blue arrows, resp.), and (iii) both cooperatively (i.e., the shifts induced by the near and far eCRFs are in the same direction, amplifying the perceptual shift). See Table S2 for a version of this table populated with representative psychophysics studies for each individual case.

# Supplementary Online Material

623

## 624 **Opponent surrounds explain diversity of contextual** 625 **phenomena across visual modalities**

626 David A. Mély<sup>1,3</sup> & Thomas Serre<sup>1,2,3</sup>

627 <sup>1</sup>*Department of Cognitive, Linguistic & Psychological Sciences*

628 <sup>2</sup>*Brown Institute for Brain Science*

629 <sup>3</sup>*Brown University, Providence, RI 02912, USA.*

630 The Supplementary Online Material includes supplementary methods, experiments and discussion  
631 with tables and figures.

### 632 **1 Supplementary Methods**

633 For all experiments, sensory variables could be read out from population responses from any  
634 column whose CRF overlapped with the center stimulus (in practice we used the column located  
635 exactly in the middle of the center stimulus). Stimuli were scaled such that the surround overlapped  
636 with both the near excitatory and the far inhibitory eCRF.

637 *Orientation tilt*

638 We modeled orientation-selective units with bell-shaped tuning curves which tile the visual field.  
639 The tuning bandwidth was  $23^\circ$ . The shift in orientation perceived by the model was computed  
640 as the difference between the orientation decoded from the center of a stimulus, with and without  
641 presenting the surround stimulus.

642 *Motion direction tilt*

643 We modeled populations of motion direction selective units with bell-shape tuning. We used the  
644 same stimuli as in the original study by Kim and Wilson (1997), holding the direction of motion  
645 of the center stimulus fixed, and varying the direction of motion of the surround stimulus between  
646  $0^\circ$  and  $180^\circ$ . The tuning bandwidth was set to  $72^\circ$ . The shift in motion direction perceived by  
647 the model was computed as the difference between the direction decoded from the center of the  
648 stimulus, with and without presenting the surround stimulus.

649 *Hue tilt*

650 We considered population units based on the V1 opponent color channels described in (Zhang  
651 et al., 2012). Isoluminant center-surround stimuli with uniformly and independently sampled hues  
652 were used for both the center and surround as done in the original study by Klauke and Wachtler  
653 (2015). The shift in hue perceived by the model was computed as the difference between the hue  
654 decoded from the center of the stimulus, with and without presenting the surround stimulus.

655 Instead of using a population vector model as done for the other modalities, we used cross-validated

656 ridge regression to decode the sine and cosine of hue. For training and testing, we used a dataset of  
657 144 colored patches uniformly sampled in the hue domain ( $n = 36$ ), presented on an approximately  
658 isoluminant achromatic background, at two saturation levels (0.20, 0.50) and two value levels  
659 (0.60, 1.0). Varying those levels or adding more of them had no significant effect either on the  
660 results or the accuracy of the decoding. We shuffled all the stimuli and randomly excluded 20%  
661 of them from the training procedure for testing; accuracy on the resulting test set was always near  
662 perfect (well above  $R^2 = 0.95$ .)

### 663 *Depth induction*

664 We modeled this experiment using population of units tuned to binocular disparity. We decoded  
665 depth from population responses centered on the test stimulus as a function of the lateral separation  
666 between the center stimulus and the flankers. We used stimuli similar to those used in the  
667 original study by Westheimer and Levi (1987), where a center line, presented at zero disparity,  
668 was flanked by two smaller squares, presented at crossed or uncrossed disparities. The tuning  
669 bandwidth was 1.15 minutes of visual angle. To assess the strength of the illusion, we used the  
670 “induction coefficient” measure as done in (Westheimer, 1986; Westheimer and Levi, 1987): we  
671 first computed the difference between the decoded disparities for the center stimulus, with and  
672 without presenting the flanker stimuli. We then reported that difference as a percentage of the  
673 difference in disparity between center and flankers.

674 *Motion direction induction*

675 We modeled this experiment using the same stimuli and protocol as used in the original study  
676 by Murakami and Shimojo (1996). The tuning bandwidth was set to  $57^\circ$ . We measured  
677 the perceptual shift through the Point of Subjective Equality (PSE), i.e., corresponding to the  
678 coherence of the center moving dots (or center coherence) at which the model had an equal  
679 probability of detecting either direction of motion. For each stimulus size, estimates of the  
680 PSE were derived from fitted model psychometric curves, each computed across 25 trials for  
681 21 regularly spaced values between  $-100\%$  and  $100\%$ . Similarly, we also computed the Point  
682 of Maximum Variance (PMV) across trials, i.e., the center coherence at which the population  
683 code for the center motion direction was maximally flat (ambiguous). Following the prescription  
684 of Ma et al. (2006), the variance of the population response curve was interpreted as a measure  
685 of uncertainty in the value decoded under the ideal neural observer model; i.e., narrow-band  
686 population response curves resulted in estimates of the sensory variable associated with high  
687 confidence levels, whereas wide-band or flat population response curves resulted in low-confidence  
688 estimates. The PMV reliably matched the PSE, so we averaged the PMV and PSE together and  
689 across trials to yield the final estimate of the PSE for each stimulus size.

690 *Enhanced color induction*

691 We modeled this experiment using the same stimuli as used in the original study (Shevell and  
692 Monnier, 2005). We used the  $b^*$  axis of the CIE  $L^*a^*b^*$  color space to measure perceptual shifts  
693 along the blue-yellow opponent axis, defined as the difference between the  $b^*$  coordinate decoded

694 from the test stimulus and the true  $b^*$  value of the test stimulus.

## 695 **2 Supplementary Experiments**

### 696 *Cross-orientation normalization*

697 Inhibition in the CRF has long been documented through cross-orientation normalization (Heeger,  
698 1993; Carandini and Heeger, 1994), where the CRF response to a primary oriented grating  
699 may be suppressed by superimposing an orthogonal masking grating. A widely-known model  
700 of single-cell and population responses alike (divisive normalization, reviewed in Carandini  
701 and Heeger, 2012), has been successful at capturing the remarkable contrast dependency of  
702 cross-orientation normalization. When the two gratings are presented at approximately equal  
703 contrasts, the population response to the resulting plaid is best described as the average of the  
704 population responses to either gratings if they were presented alone (the summation regime).  
705 However, if the two gratings contrasts are disparate enough, the resulting population response  
706 looks like the response to the grating with the strongest contrast, as if the weaker-contrast grating  
707 was not presented (the winner-take-all regime).

708 Our model, which includes untuned inhibition within the CRF, captures the balance between the  
709 summation vs. the winner-take-all regimes and its dependence on both gratings' contrasts. We  
710 digitally extracted electrophysiology data from Figure 4A, 4E and 4G from (Busse et al., 2009) in  
711 Figure S2 and show the corresponding model data in Figure S3.

712 *Feature-matched surround suppression*

713 Center-surround stimuli are widely used in neurophysiology to probe the mechanisms underlying  
714 surround suppression (Hubel and Wiesel, 1968; DeAngelis et al., 1994; Weliky et al., 1995; Petrov  
715 et al., 2005; Ozeki et al., 2009). In many studies, a cell is driven by a center grating at the  
716 optimal orientation; presenting a surround grating with the same orientation (or extending the  
717 center stimulus beyond the CRF and into the eCRF) suppresses the CRF response.

718 However, since the center grating was chosen to be optimal for the recorded cell, it remains unclear  
719 whether the surround grating is maximally suppressive because it matches the cell's orientation  
720 preference, or because it matches the orientation of the center grating. Two studies in V1 (Trott  
721 and Born, 2015; Shushruth et al., 2012) support the latter hypothesis, i.e., surround suppression is  
722 strongest when the surround stimulus matches the center in orientation, irrespective of the cell's  
723 orientation preference. This excludes simpler models where a center cell is simply selectively  
724 targeted by tuned inhibitory connections from its surround (Schwabe et al., 2006).

725 Similarly to Shushruth et al. (2012), we have found that this phenomenon is best accounted for by  
726 weakly-tuned recurrent excitation within an (orientation) hypercolumn. Such excitation places the  
727 hypercolumn in a highly recurrent regime where a cell may be driven by a stimulus at a sub-optimal  
728 orientation. Then, the most suppressive surround stimulus matches the center orientation because  
729 it withdraws most of the cell's input.

730 Our model explains neurophysiology data (digitally extracted from Figure 1B of Trott and Born  
731 (2015), curve labelled as “surround tuning”) for such feature-matched surround suppression (see

732 Figure S3). It also explains another experiment where the presentation of a surround grating  
733 selectively suppresses population responses to the identically-oriented component of a two-grating  
734 plaid presented in the CRF, in effect “cancelling” the contribution of the plaid component matched  
735 to the surround (see Figure S5).

### 736 *The necessary mechanisms of the contextual model*

737 Our model successfully recreated human perceptual and primate neurophysiological responses  
738 across a variety of contextual phenomena. This suggests that the model’s full suite of mechanisms  
739 are sufficient for a general explanation of contextual processing, but are they all necessary? We  
740 investigated this with a large-scale lesion screening procedure that identified the minimal version of  
741 the model. This procedure involved lesioning a mechanism in the model by setting its weights to 0,  
742 and then optimizing the model’s remaining free parameters for explaining contextual phenomena.  
743 We hereafter refer to this as the lesioning procedure. We used the lesioning procedure to separately  
744 measure the necessity of each of the model’s critical components. In total, we tested: (1) the  
745 necessity of each CRF and eCRF mechanism, (2) the validity of our assumptions about their  
746 spatial configurations, (3) the importance of asymmetric facilitatory vs. suppressive contributions  
747 for explaining contextual phenomena, and (4) the sensitivity of the model to a range of stimulus  
748 tuning properties and relaxed constraints on its patterns of connectivity. We performed this analysis  
749 on all contextual phenomena explored in the main and supplemental text except for Figure 6, which  
750 was excluded because its computational complexity rendered it intractable for this procedure.

751 The lesioning procedure was performed separately on each component of the model. We also

752 applied the free parameter optimization routine to the full model without lesioning any of its  
 753 components. This supported identification of what parameters most accurately explained each  
 754 of the contextual phenomena. This procedure is equivalent to solving:

$$C_{lesion} = \arg \max_{\forall j \in J} f(\rho_j) \quad (10)$$

755 Where  $C_{lesion}$  is the selected model from the optimization procedure over  $J = 1,000$  iterations  
 756 for a specific model configuration (e.g., a model with a lesion to its far eCRF). This optimization  
 757 is performed over hand-tuned model parameters  $\alpha, \beta, \mu, \nu, \gamma, \delta$  in the space of all possible  
 758 combinations of these parameters across the perceptual phenomena discussed in the main text.  $\rho$   
 759 is a vector of correlation scores for a selected model configuration describing the quality of its  
 760 simulations for each perceptual phenomenon.  $f$  is a monotonic function of  $\rho$  as follows:

$$f(\rho) = \mu(\rho_j) / \sigma(\rho_j) \quad (11)$$

761 Where  $\mu$  calculates the mean correlation and  $\sigma$  calculates the standard deviation of correlations  
 762 across perceptual phenomena for a sampled model configuration. In effect,  $f$  penalizes models  
 763 by  $\sigma$  for overfitting on a subset of perceptual phenomena, making it a better method than  
 764 simply maximizing correlation for selecting model configurations that are generally successful  
 765 at simulating perception. Pearson correlation is desirable for this optimization procedure because  
 766 it does not penalize irrelevant differences between model data and perceptual data, such as shifts  
 767 of the mean and rescaling.

768 Each set of parameters applied to a model configuration was selected by searching over an  
769 exponentially spaced grid around each of the hand-tuned model parameters listed above. We define  
770 this sampling procedure  $S$ , which was applied on every iteration to each of the hand-tuned model  
771 parameters  $p$  discussed in the main text:

$$S(p) = |p + \text{unif}[-1, 1]^{\text{unif}[-2, 2]}| \quad (12)$$

772 where *unif* denotes uniform sampling in the specified range. The parameters that maximize  $C_{\text{lesion}}$   
773 yield optimal performance (as measured by Pearson correlation) for a lesioned version of the model  
774 in recreating observers' responses across contextual phenomenon.

775 While we found all model mechanisms to be necessary, we did find that the full version of the  
776 model was relatively tolerant to parameter perturbations caused by this optimization procedure.  
777 See Figure S9 for histograms of these scores. Plotting performance of the full model across these  
778 parameters on any particular problem revealed that the majority of them yielded simulations that  
779 were qualitatively similar to the perceptual data. Variations in performance potentially provide  
780 insight into individual differences in perception of contextual phenomena. See Figure S13 for an  
781 example of the qualitative variability we observed for a representative phenomenon.

782 A combination of quantitative and qualitative evidence demonstrates that our model performs best  
783 when it contains the full suite of CRF and eCRF mechanisms. Lesioning either of the eCRF  
784 mechanisms (near excitation or far inhibition) diminishes the model's ability to explain observers'  
785 behavior for several phenomena (Figure S8). Lesioning the near eCRF tuned excitation degrades

786 the ability of the model to explain phenomena demonstrated in Figure 3C, Figure 5, Figure 7, and  
787 Figure S5. Lesioning the far eCRF tuned inhibition diminishes the model explanatory power of  
788 Figure 3C (see also Figure S14), Figure 4, Figure 7, and Figure S5. Indeed, the full model was  
789 significantly better at capturing these phenomena than either lesioned version, as measured with  
790 2-tailed t-tests comparing the maximum correlations accrued by each model across each of the  
791 eight phenomena included in the lesion optimization (Full model vs. near eCRF excitatory lesion:  
792  $t(7) = 5.102, p = 0.001$ ; Full model vs. far eCRF inhibitory lesion:  $t(7) = 3.911, p = 0.006$ ).

793 Lesions to either of the model's CRF mechanisms reveal that they were less important for  
794 explaining the high-level contextual phenomena discussed in the main manuscript than its eCRF  
795 mechanisms (Figure S7). A closer inspection of lower level phenomena, however, revealed their  
796 importance for explaining contrast-dependence tuning in V1. We measured responses from the  
797 full optimized model and versions of the model with lesioned CRF mechanisms to a stimulus that  
798 varied in size and contrast. Lesioning either the weakly tuned excitation or untuned inhibition CRF  
799 mechanisms qualitatively harmed the model's ability to discriminate between stimuli of different  
800 contrasts (Figure S1). For the model with lesioned weakly tuned excitation, this is immediately  
801 apparent: it is qualitatively worse than the full model at discriminating stimulus contrast until the  
802 size of the stimulus extends into the near eCRF. For the model with lesioned untuned inhibition,  
803 the opposite phenomenon is observed: contrast discrimination fails as soon as the stimulus extends  
804 into the inhibitory far eCRF.

805 Lesion optimization also highlighted the importance of our key assumption of asymmetry between  
806 excitation versus inhibition for explaining contextual phenomena. Excitation in the model is purely

807 additive and only depends on pre-synaptic activity. But inhibition depends on both pre- and  
808 post-synaptic activity and results in a combination of subtractive and divisive effects (Carandini  
809 and Heeger, 2012). A model with lesioned presynaptic shunting inhibition was significantly worse  
810 at explaining the contextual phenomena than the full model ( $t(7) = 2.631$ ,  $p = 0.034$ ; CRF  $\alpha$  and  
811 far eCRF  $\beta$ , Figure S8).

812 This lesion-screening framework allowed us to measure the necessity of having spatially separate  
813 versus overlapping CRF and near eCRF mechanisms. We created a version of the model in which  
814 the influence of CRF and near eCRF mechanisms were averaged together. This configuration  
815 yielded significantly worse than the full model at explaining the contextual phenomena ( $t(7) =$   
816  $2.908$ ,  $p = 0.023$ ; Spatially overlapping CRF and near eCRF, Figure S8 ). The impact of having  
817 spatially distinct CRF and near eCRF mechanisms was most apparent on Figure 4A-B, Figure 5  
818 and Figure 7.

819 We also investigated the necessity of separate near and far eCRF regions for explaining these  
820 perceptual phenomena. We optimized a model with completely overlapping near and far eCRFs –  
821 extending from the proximal point near eCRF to the distal point of the far eCRF. The spatially  
822 separate eCRFs of the full model yielded significantly better performance than this spatially  
823 overlapping version at explaining contextual phenomena ( $t(7) = 2.552$ ,  $p = 0.038$ ; Spatially  
824 overlapping near eCRF and far eCRF, Figure S8).

825 Having found evidence that each of the contextual model's mechanisms are necessary to explain  
826 the full array of contextual phenomena, we explored its robustness to variations in the shape

827 of the model unit tuning curves. We did this over 1,000 iterations by resampling the tuning  
828 curve properties for each contextual phenomenon with the sampler  $S$ . For disparity (Cumming  
829 and Parker, 1997), motion direction (Albright et al., 1984) and orientation (Ringach et al., 1997)  
830 this involved resampling tuning curve bandwidth with the sampler  $S$ . For color opponent tuning  
831 phenomena, for which we did not assume a bell-shape tuning curve, this involved resampling the  
832 response threshold from a uniform distribution in  $[-0.5, 0.5]$ . On every iteration of the procedure,  
833 we simulated each contextual phenomenon with the optimized full model after resampling its  
834 tuning properties. This approach revealed that, with the exception of Figure 4, the full optimized  
835 model was robust to a wide range of tuning properties. See Figure S11 for histograms of the  
836 model’s ability to explain contextual phenomena as these parameters were varied.

837 The contextual model has strictly “tuned” connections from the eCRF onto the CRF, with a  
838 one-to-one mapping between computational units preferring the same stimulus features. We tested  
839 how important this constraint is to the model’s performance by relaxing this strict one-to-one  
840 mapping into a weakly tuned eCRF-CRF mapping (Figure S10). eCRF units in the orientation,  
841 motion, and disparity domains had normally distributed connectivity, centered at a target unit  
842 tuning preference  $\theta_k$  with standard deviation  $\zeta$ . Because we consider color-opponent eCRF units  
843 with monotonic tuning curves in the color domain, assumptions of normality are inappropriate.  
844 For these phenomena we instead set  $w_{kk} = (\zeta\sqrt{2\pi})^{-1}$  and  $w_{jk} = \text{const.}$  (when  $j \neq k$ ; under the  
845 constraint that the weights sum up to  $r1$ ). In each case the lesion screening procedure was used  
846 to search the extent to which eCRF unit connectivity could be weakened without destroying the  
847 model’s ability to simulate contextual phenomena. Model performance was recorded over 1,000

848 iterations while the standard deviation  $\zeta$  of these connectivity schemes was randomized with the  
849 sampler  $S$  (see 12 for details on the sampling procedure). For each contextual phenomenon this  
850 procedure yielded many successful simulations, and at times outperformed strict tuning. We expect  
851 that additional work on incorporating more anatomically plausible connectivity into the full model  
852 will yield even better performance than we report here.

853 Taken together, our large-scale lesion screening procedure indicates that the mechanisms in our  
854 full model are not only sufficient for explaining contextual phenomena, but also necessary.

### 855 **3 Supplementary Discussion**

#### 856 *Neurophysiology evidence for shifts in single-cell tuning curves*

857 A key prediction of the model is that shifts in population responses may underlie perceptual shifts  
858 in phenomena as in the two regimes of the tilt effect. Several electrophysiology studies have  
859 provided evidence for shifts in single-unit tuning curves that matched behavioral data on perceptual  
860 shifts. In the motion domain, Li et al. (1999) recorded from speed-tuned neurons in V1 and found  
861 that the presence of a slower (faster) moving stimulus in a neuron's eCRF shifts its preferred speed  
862 towards faster (slower) values. In the color domain, shifts in tuning curves that are consistent with  
863 color contrast were found both in V1 (Wachtler et al., 2003) and V4 (Kusunoki et al., 2006). In  
864 two of these studies (Li et al., 1999; Kusunoki et al., 2006), single-unit recordings were shown  
865 to be consistent with the animals' behavioral responses. In the disparity domain, Thomas et al.  
866 (2002) recorded from single cells in monkey V2 and found that presenting stimuli with different  
867 disparities in their CRF and eCRF resulted in tuning curve shifts towards the disparity presented

868 in the eCRF. This behavior is consistent with the computation of relative disparity. Interestingly,  
869 this yields a rather counter-intuitive model prediction for the associated perception: single-cell  
870 tuning curve shifts *towards* the eCRF disparity translate, at the population level, to a shift of the  
871 estimated disparity in the CRF *away* from that of the eCRF. That is, apparent attraction at the  
872 level of individual tuning curves effectively corresponds to repulsion at the perceptual level; this  
873 paradox emphasizes the necessity for a population-level analysis to fully appreciate the perceptual  
874 effect of extra-classical modulation.

875 *Evolution of perceptual shifts as a function of population response curve bandwidths*

876 In the orientation and motion direction tilt effects, our model suggests that repulsive shifts happen  
877 when the eCRF suppresses one side of the CRF population response (“pushing” it away from its  
878 own mode) and that the strength of such a push grows with the overlap between these populations  
879 (because of the inhibition increasing with post-synaptic activity). Consistently, Goddard et al.  
880 (2008) have found that in the orientation tilt effect, using stimuli embedded in noise, broadening  
881 the orientation power of the center stimulus (yielding broader-band center population response  
882 curves) results in larger repulsive shifts than with narrower-band stimuli (e.g., the gratings used  
883 in (O’Toole and Wenderoth, 1977)).

884 Figure S12 shows a prediction of the model regarding how the largest achievable repulsive shift  
885 should evolve as a function of the bandwidths of the center and the surround populations. For  
886 example, large effects are predicted when the center bandwidth is much larger than the surround  
887 bandwidth (up to a factor of 2), presumably because high overlap between population responses

888 from the CRF and the eCRF can then be achieved while keeping them offset from each other in the  
889 orientation domain (perfect alignment would result in equal inhibition on either side of the CRF  
890 population response curve, yielding no shift).

891 Conversely, our model predicts that the attractive regime of the tilt effect can only exist when CRF  
892 and eCRF populations do not overlap significantly. By extension, this implies that the attractive  
893 regime weakens, or even disappears altogether, when typical population bandwidths are so large  
894 that there is no “room” left in the tuning domain for the CRF and eCRF populations to exist  
895 without overlapping. This is consistent with the observation that the attractive regime in the  
896 motion direction tilt effect is much less pronounced than in the orientation tilt effect (Figure 3), as  
897 neurophysiology data from V1 suggests that direction-tuned cells have a higher tuning bandwidth  
898 than orientation-tuned cells (Ringach et al., 2002; Albright et al., 1984).

#### 899 *Surround inhibition explains perceptual repulsion*

900 To explain color induction, Klauke and Wachtler (2015) have proposed a phenomenological model  
901 based on weakly-tuned surround suppression, which reduces the gain of individual bell-shaped,  
902 hue tuning curves in the center. This results in center population shifts that are qualitatively  
903 consistent with their behavioral data. Our model thus provides computational evidence for their  
904 population-level explanation, using a similar form of surround suppression with sharper tuning.  
905 Interestingly, we show that neural populations with explicit bell-shaped hue tuning are actually  
906 not necessary and perceptual shifts consistent with psychophysical data can be accounted for with  
907 neural populations tuned to cardinal color opponency. Additionally, our experiments involving

908 visual modalities with bell-shaped tuning curves (e.g., orientation, motion direction tilt effects) also  
909 confirmed their intuition that tuning bandwidth controls the maximal amplitude of the perceptual  
910 shift.

911 Goddard et al. (2008) also explains the orientation tilt effect with a computational model based  
912 on tuned surround inhibition. Their model also incorporates a divisive normalization term but  
913 lacks a realistic model of recurrent connections. In order to explain the emergence of an attractive  
914 regime at dissimilar orientations in the center and the surround, these authors modeled the surround  
915 orientation tuning curve as a “Mexican hat”, that has a positive peak around the center orientation,  
916 and becomes negative for values further away. This is superficially similar to the balance between  
917 the influences of the near and far eCRFs in our model. However, such a model assumes that  
918 excitatory and inhibitory influences are evenly distributed across the surround, and thus would  
919 fail to explain the dependence of attraction vs. repulsion on stimulus size in our second set of  
920 illusions. Interestingly, the authors allow the relative shape of the surround tuning curve (i.e.,  
921 the excitation-inhibition balance), as well as the overall strength of surround modulation, to vary  
922 across experiments to fit experimental data. This suggests that carefully allowing the analogues  
923 of these parameters in our model to vary across visual modalities may result in closer quantitative  
924 fits.

#### 925 *Enhanced color shift effect: beyond S-cone contrast*

926 Monnier and Shevell (2003); Shevell and Monnier (2005) modeled enhanced color shifts through  
927 an antagonist CRF organization corresponds to an S-ON center (i.e., the neuron response increases

928 with short cone activation in its center) and S-OFF surround (i.e., the neuron response decreases  
929 with short cone activation in its surround). Their explanation is consistent with ours: at the optimal  
930 spatial frequency (which exists because spatial frequency controls the overlap of the different parts  
931 of the stimulus with the different sub-regions of the CRF or the eCRF), the two colors present in  
932 the contextual stimulus create cooperative shifts in S cone activation. However, such a model does  
933 not predict enhanced perceptual shifts using patterned contextual stimuli that have zero contrast in  
934 S-cone activation; the authors report little perceptual shifts along L and M cone activations, which  
935 could be due to the specific stimuli they used. We rendered a simple equivalent of their stimuli  
936 with contextual colors that should elicit little to no contrast in S-cone activations (Figure S6E-H).  
937 In other words, CRFs based on S-cone activations should respond uniformly across the surround  
938 stimulus of this modified illusion; yet, the shift seems just as vivid (see caption for details) as with  
939 the original stimuli (Shevell and Monnier, 2005). Interestingly, the model reported in (Shevell and  
940 Monnier, 2005) to best account for behavior after parameter fitting results in a surround up to 5  
941 times wider than the center, and a center almost 2 times as responsive to visual stimulation as the  
942 surround. This seems to agree better with an eCRF model like ours, which has an inhibitory far  
943 eCRF that is (relatively) much wider than the OFF region of traditional center-surround CRFs.

	<b>parameter description</b>	<b>value</b>
$\eta$	Input time constant	6.00
$\varepsilon$	Input gain	0.50
$\xi$	Afferent strength	4.50
$\tau$	Output time constant	6.00
$\sigma$	Output gain	0.50
$\alpha$	Untuned suppression strength (divisive)	1.00
$\mu$	Untuned suppression strength (subtractive)	1.00
$\beta$	Tuned suppression strength (divisive)	3.00
$\nu$	Tuned suppression strength (subtractive)	0.30
$\gamma$	Tuned facilitation strength	1.00
$\delta$	Untuned facilitation strength	1.00
$\varsigma$	Standard deviation of tuned facilitation weights	0.15
$R[\mathbb{N}_{x,y}]$	Radius of near eCRF of $(x,y)$ (in number of hypercolumns)	9
$R[\mathbb{F}_{x,y}]$	Radius of far eCRF of $(x,y)$ (in number of hypercolumns)	29

Table S1. *Model parameters*. Hypercolumns are organized in a regularly-spaced, square grid.

Radii for near and far eCRFs are then defined in number of hypercolumns on that grid.

	<b>Separate</b>	<b>Competition</b>	<b>Cooperation</b>
<b>Shifts</b>	Goddard et al. (2008) (o.), Westheimer (1986); Westheimer and Levi (1987) (d.)	Gibson and Radner (1937); O’Toole and Wenderoth (1977) (o.), Marshak and Sekuler (1979) (m.), Loomis and Nakayama (1973); Norman et al. (1996); Baker and Graf (2010) (s.)	Monnier and Shevell (2003); Shevell and Monnier (2005) (c.)
<b>Bumps/Notches</b>	Goddard et al. (2008) (o.), Murakami and Shimojo (1993; 1996) (m.)	Smith et al. (2001); Klauke and Wachtler (2015) (c.)	Anstis (2006) (b., c.)

Table S2. *How existing psychophysics studies on contextual phenomena across modalities fit in the proposed taxonomy.* See Figure 8 in the main text for a companion figure. **Abbreviations:** b: brightness, c: color, d: depth, m: motion direction, o: orientation, s: motion speed.

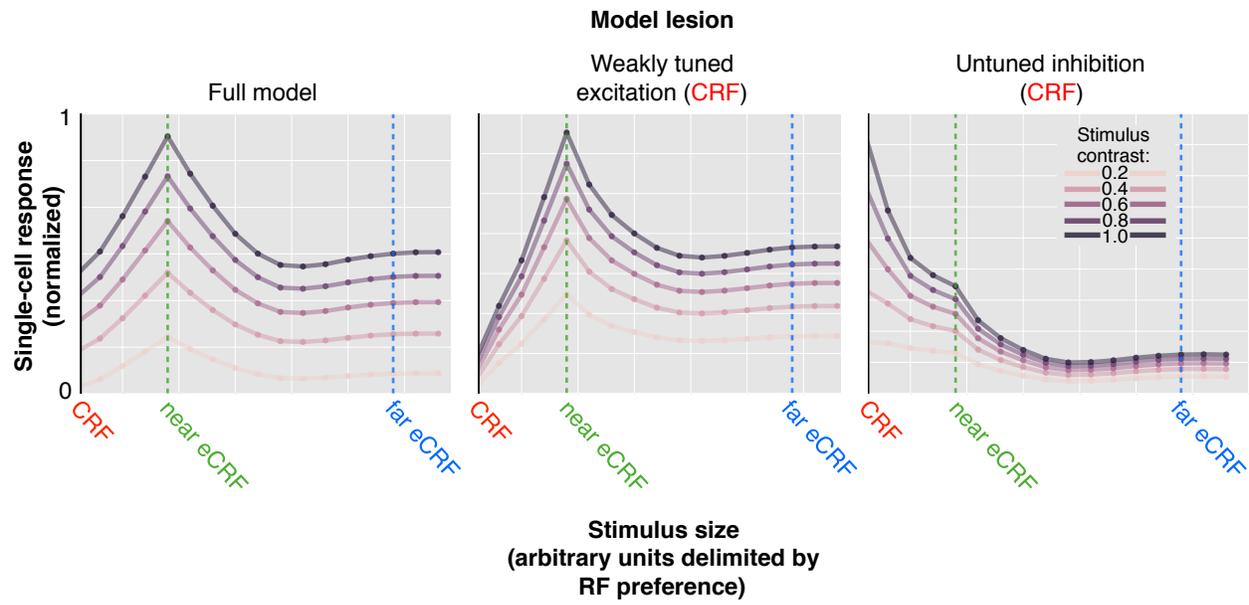


Figure S1. ***Both intra-columnar recurrent excitation and inhibition are necessary for explaining neural contrast responses.*** Model response is plotted as a function of the size (x-axis) and contrast (line color) of a stimulus. Vertical dashed lines show the spatial extent of the CRF (red) and near excitatory (green) and far inhibitory (blue) eCRFs. All remaining model parameters were optimized for the contextual phenomena in the main manuscript after selectively lesioning either of the two mechanisms. Only the full model exhibits expected tuning properties: the CRF response discriminates stimuli and reaches its peak response once stimuli expand into the excitatory near eCRF receptive field.

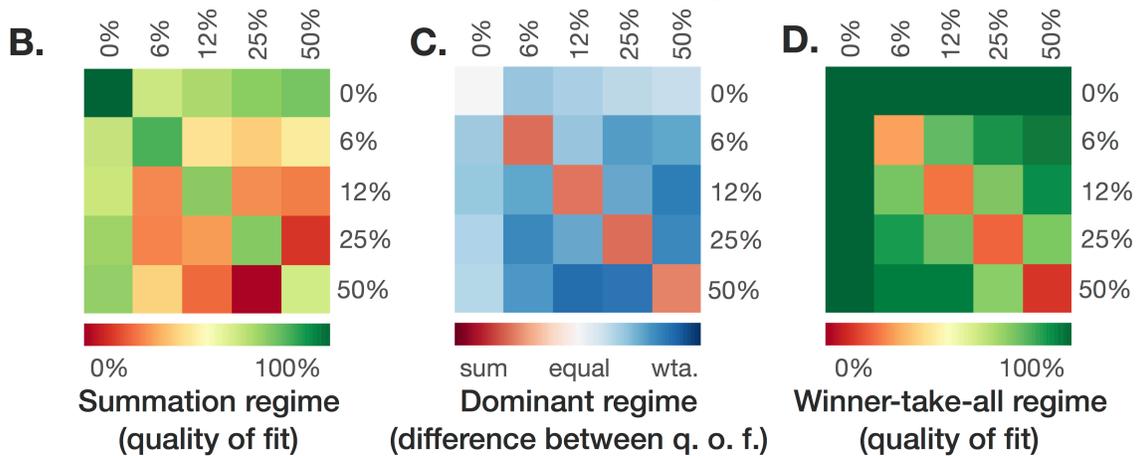
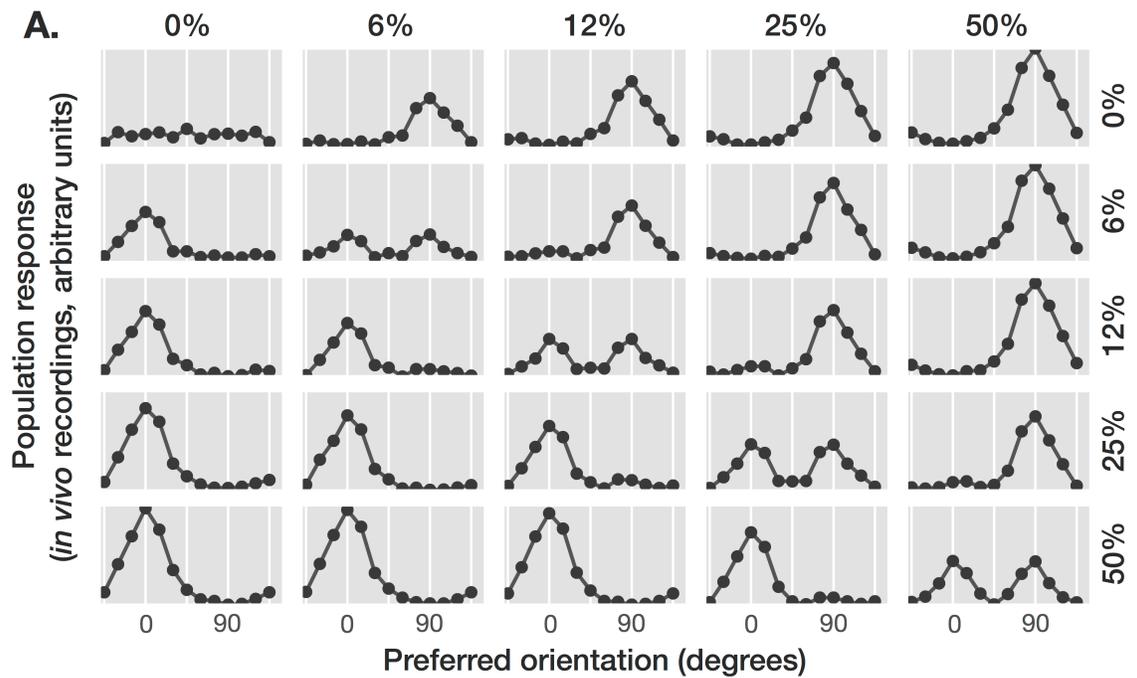


Figure S2. *Balance between summation and winner-take-all regimes during the presentation of cross-oriented gratings.* **Upper plots:** cat V1 electrophysiology data digitally extracted from Figure 4A in (Busse et al., 2009). Population recordings of orientation-tuned cells in response to the presentation of orthogonally-oriented gratings at variable contrast levels (rows indicate the contrast of the  $0^\circ$  grating and columns indicate the contrast of the  $90^\circ$  grating). When the contrasts of the two gratings are approximately equal, the population response is best described as the sum of the population responses to either grating presented in isolation. When they differ markedly, the population response is best described as the population response to the grating with the strongest contrast presented in isolation. This latter regime is also known as a winner-take-all regime as it seems that the response to the stronger stimulus predominates and suppresses the response to the weaker one. **Lower plots:** data reproduced from Figure 4E and Figure 4G in (Busse et al., 2009). Quality of fit of either a summation model or a winner-take-all model as a function of the contrasts of either component grating. As before, the relative strengths of either component grating controls

945 the balance between summation and winner-take-all.

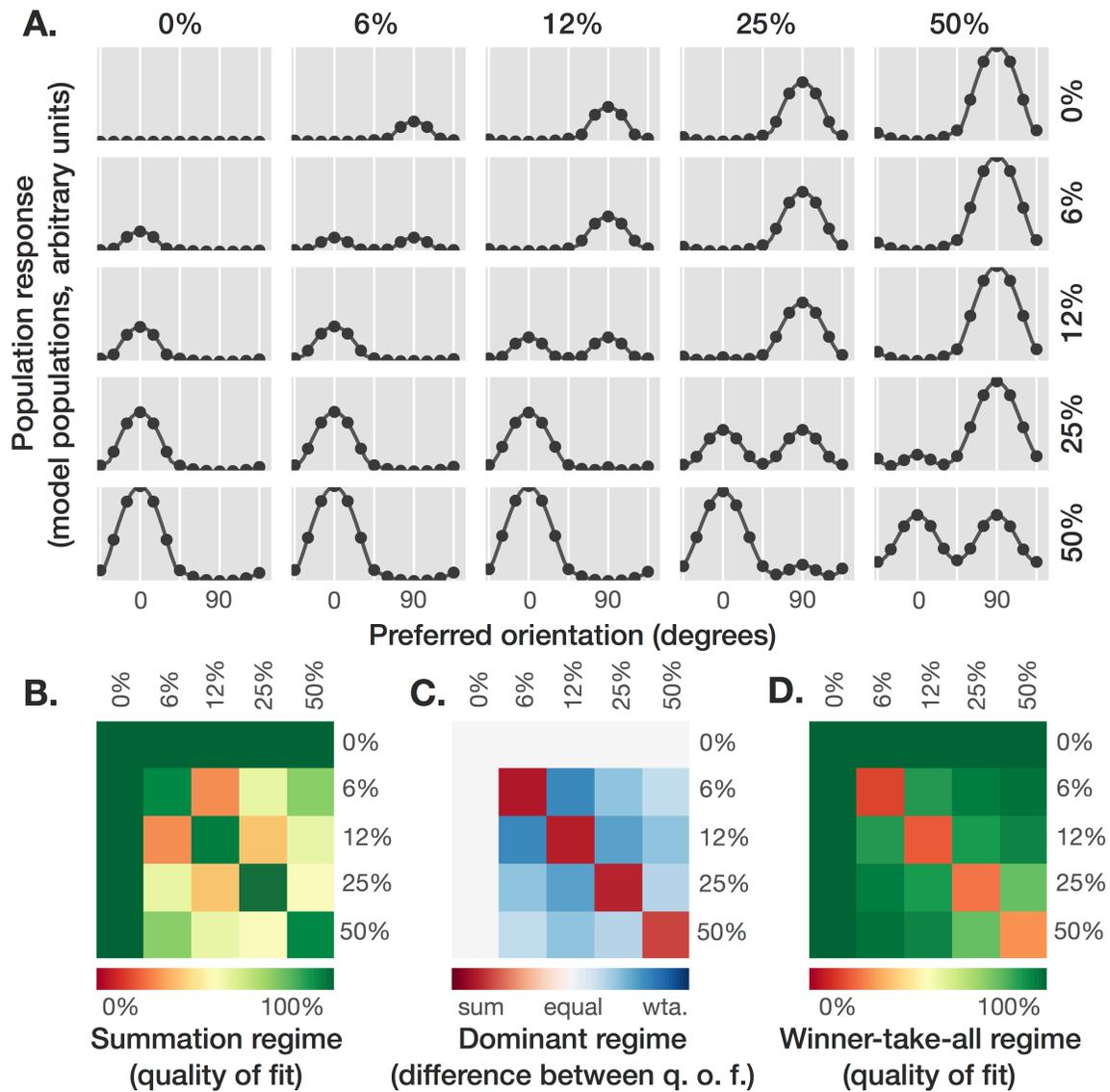


Figure S3. *Balance between summation and winner-take-all regimes during the presentation of cross-oriented gratings in the model.* Using populations of orientation-tuned cells, our model

946 reproduces the neurophysiology data presented in Figure S2.

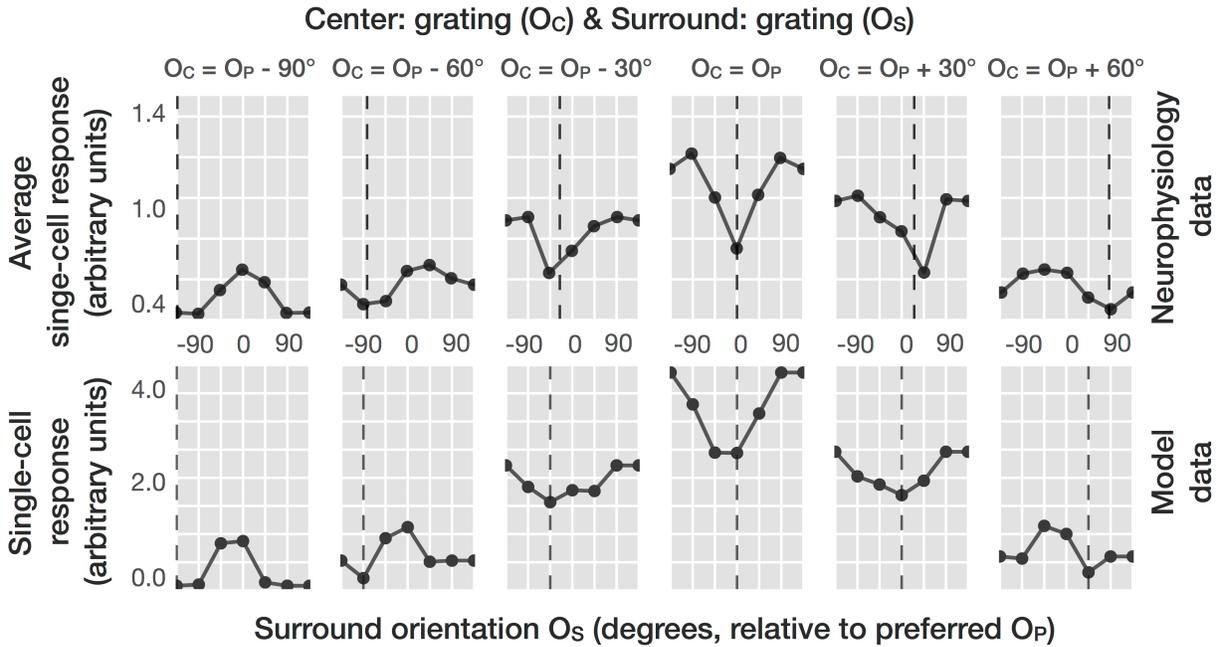


Figure S4. *Feature-specific surround suppression: the most suppressive surround stimulus matches that of the center stimulus.* **Top row:** monkey V1 electrophysiology data digitally extracted from Figure 1B (curve labelled as “surround tuning” in Trott and Born, 2015). Single-cell recordings from orientation-tuned cells (averaged and normalized across cells such that  $0^\circ$  on the abscissas corresponds to  $O_P$ , the cell’s preferred orientation) were made using center-surround gratings. Different center grating orientations  $O_C$  were presented in the CRF. For each value of  $O_C$ , the orientation of the surround grating  $O_S$  was systematically varied (abscissas). The key result is that the most suppressive surround orientation (vertical dashed black line) always matches the orientation of the center stimulus  $O_C$ , even for non-optimal center orientations. **Bottom row:** The model exhibits a similar feature-selective suppression as found experimentally.

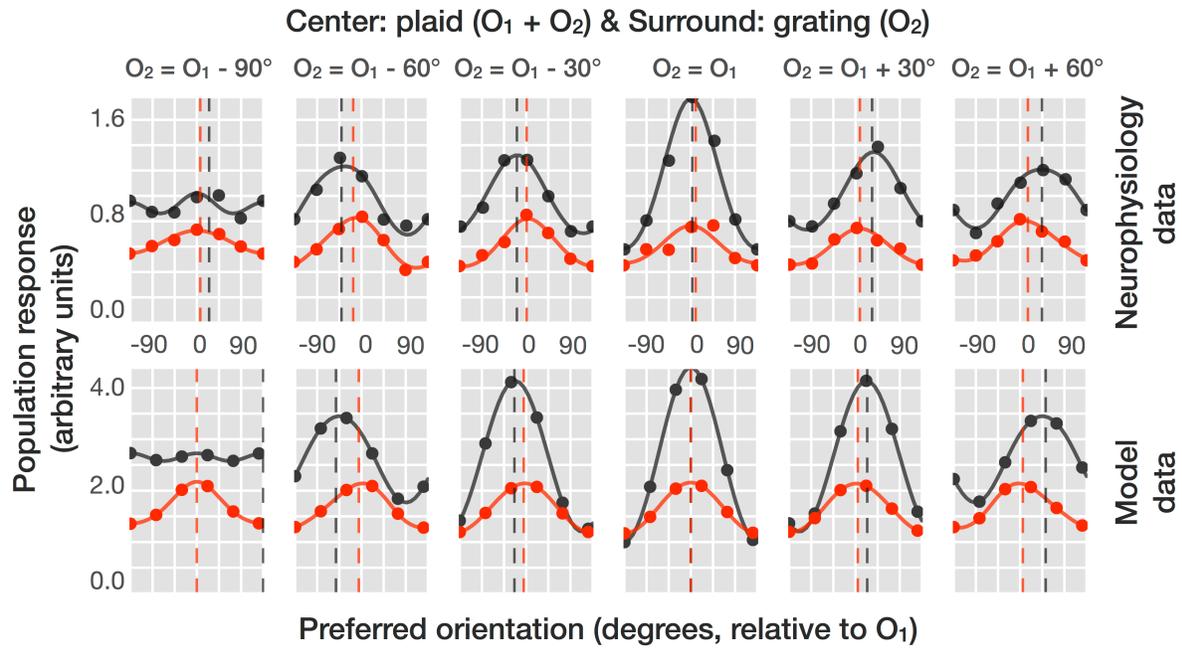


Figure S5. *Feature-specific surround suppression: the surround stimulus suppresses the iso-oriented component of a center plaid.* **Top row:** monkey V1 electrophysiology data reproduced from Figure 4 B from (Trott and Born, 2015). Population-level recordings were made from orientation-tuned cells. A plaid was first presented in the center (with components at orientations  $O_1$  and  $O_2$ ), alone. The resulting population response curves look like the averaged population responses to the presentation of either plaid component in isolation (black curves; black dashed lines placed at the orientation encoded by the corresponding populations; compare also to cross-orientation normalization, Figures S2 and S3). When a surround grating with the same orientation as one of the center plaid's components (e.g.,  $O_2$ ) was added to the surround, the corresponding plaid component (e.g.,  $O_2$ ) was selectively suppressed, resulting in a population response tuned to the other plaid component (e.g.,  $O_1$ ; red curves; red dashed lines corresponds to the orientation encoded by the corresponding populations). **Bottom row:** The model exhibits feature-selective suppression as well (the red dashed lines remain centered at  $O_1$ , the plaid component that does not match the surround orientation).

947

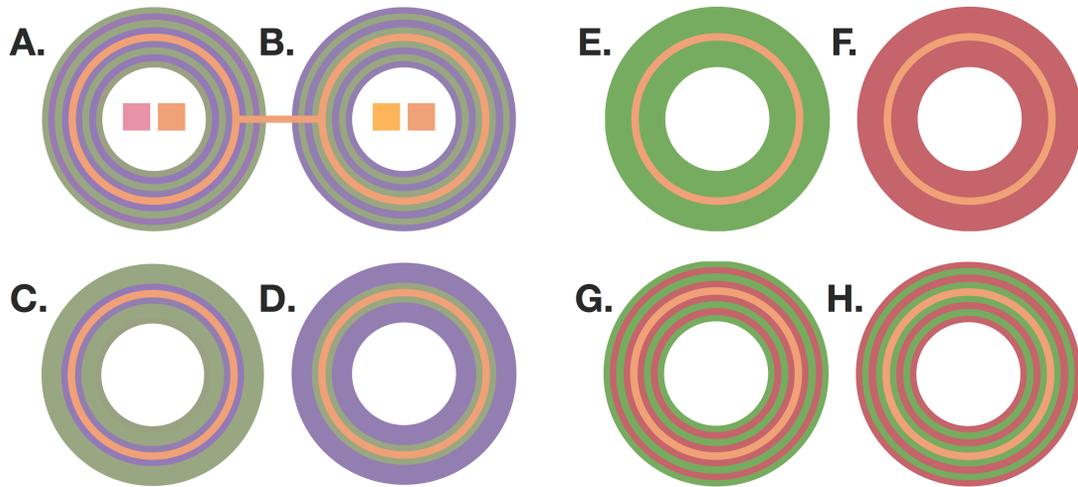


Figure S6. *Enhanced color shifts: model readout and novel predictions.* **A-B.** Stimuli were adapted from Figure 1 in (Monnier and Shevell, 2003). Though the test ring has the same color in either case (orange), patterned contextual stimuli of opposite phases result in color shifts in opposite directions. The squares in the middle of the rings represent the hue of the test ring decoded from our model (left square) and from a CRF-only model of opponent color processing for baseline comparison (Zhang et al., 2012; right square). The hue decoded from our model appears more consistent with human perception. **C-D.** Our model explains the illusion through the cooperative stimulation of the near and far eCRF by near-opposite hues (see main text and Figure 7), here purple and lime. As a result, we predict that the illusion should have the same intensity when the contextual patterned grating is replaced by two uniform regions with the right size in order to coincide with the near and far eCRFs (compare the appearance of the test ring in **C/D** to that in **A/B**). **E-H.** A previous explanation for this illusion involved a center-surround antagonism for S-cone activity (Monnier and Shevell, 2003; Shevell and Monnier, 2005). We rendered a new version of the illusion predicted by our model to be just as vivid as the original one, but whose patterned context lacks S-cone contrast (colors rendered approximately; appearances may differ on paper). Thus, whereas a S-cone antagonist model should see no difference in the appearance of the test ring in **E-H**, we predict a slight difference between **E-F**, and a very vivid difference

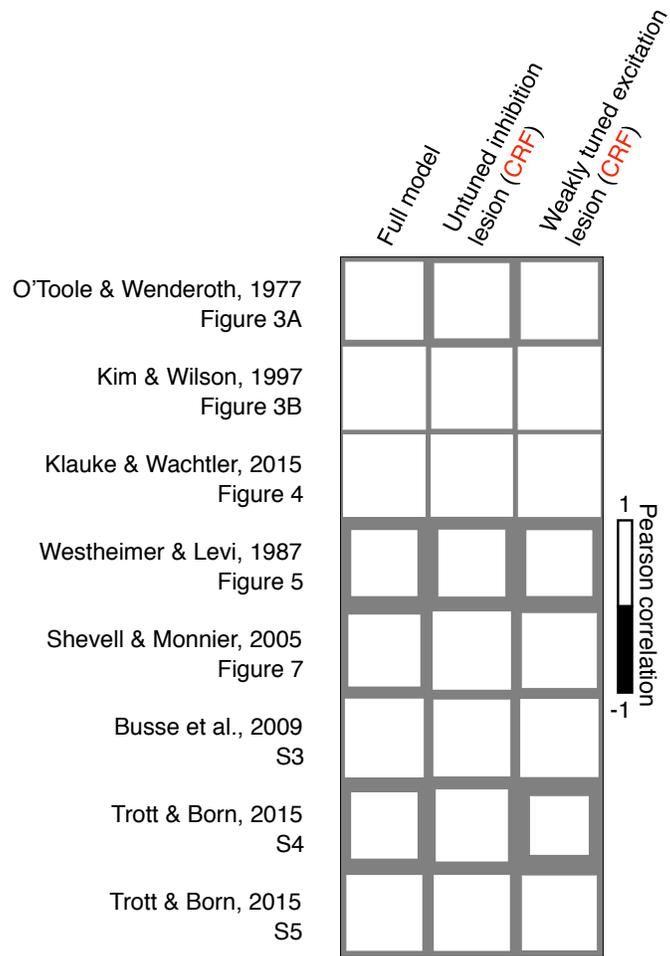


Figure S7. *Lesioning model CRF mechanisms mostly leaves its performance intact.* We measured the model's ability to explain a variety of contextual phenomena (depicted along rows) when it was intact versus when one of its CRF circuits was lesioned (depicted along columns). Correlations between each experiment's behavioral data and model simulations are plotted with a Hinton diagram, where each square's size depicts correlation magnitude and color depicts the sign of the correlation. All correlations  $> 0$ . Model simulations follow 1000 iterations of parameter optimization. Figure 6 is omitted because the complexity of that simulation made the parameter

949 search computationally intractable.

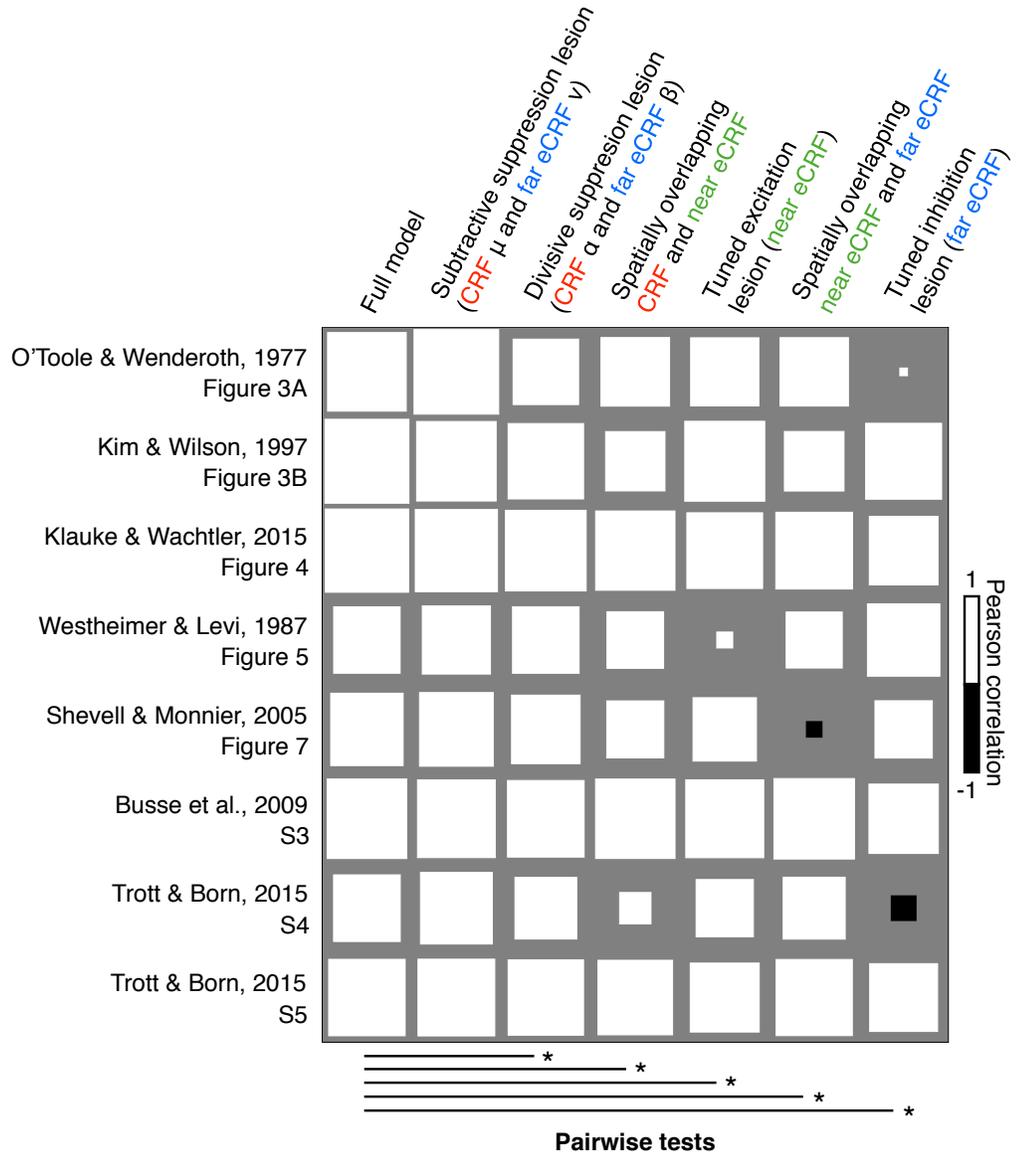


Figure S8. *Model eCRF mechanisms are necessary for its ability to explain contextual phenomena.* We measured the model's ability to explain a variety of contextual phenomena (depicted along the rows) when it was intact versus when one of its circuits was lesioned (depicted along the columns). Correlations between each experiment's behavioral data and model simulations are plotted with a Hinton diagram, where each square's size depicts correlation magnitude and color depicts the sign of the correlation. Model simulations follow 1000 iterations of parameter optimization. Figure 6 is omitted because the complexity of that simulation made the parameter search computationally intractable. Also note that we select the optimal model over all experiments which is not necessarily the optimal model for any given individual experiment. Lines and \* at the bottom identify lesioned versions of the full model with correlations to perceptual data that are significantly worse than the full model. Comparisons are made with two-tailed *t*-tests. All

950 tests shown are  $p < 0.05$ .

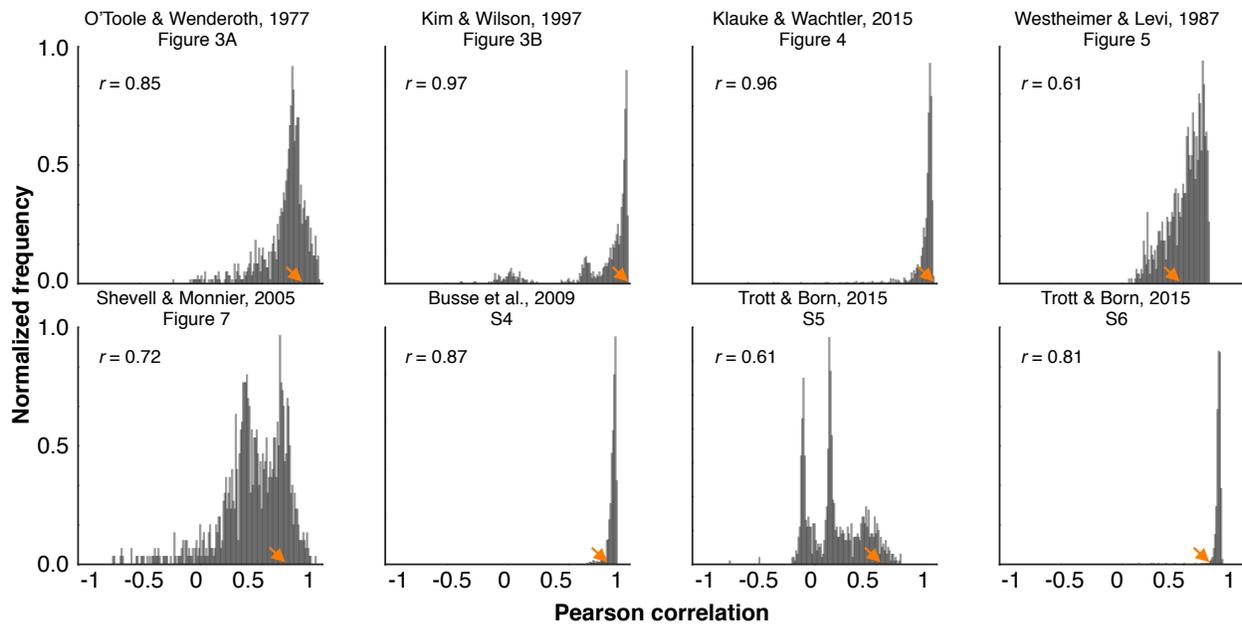


Figure S9. *Distributions of fit derived from parameter optimization of the full model.* Histograms depict accuracy of the full model in explaining each contextual phenomenon as its free parameters are randomly sampled with an exponential search grid over model parameters. Performance is measured with Pearson correlation. Correlations are derived for each of the 1000 iterations of the lesion-screening procedure. Correlations reported in the top left-hand corner of each panel denote performance of the optimized full model on each problem. These values are also highlighted with orange arrows. Figure 6 is not included because it is computationally intractable for the

951 lesion-screening.

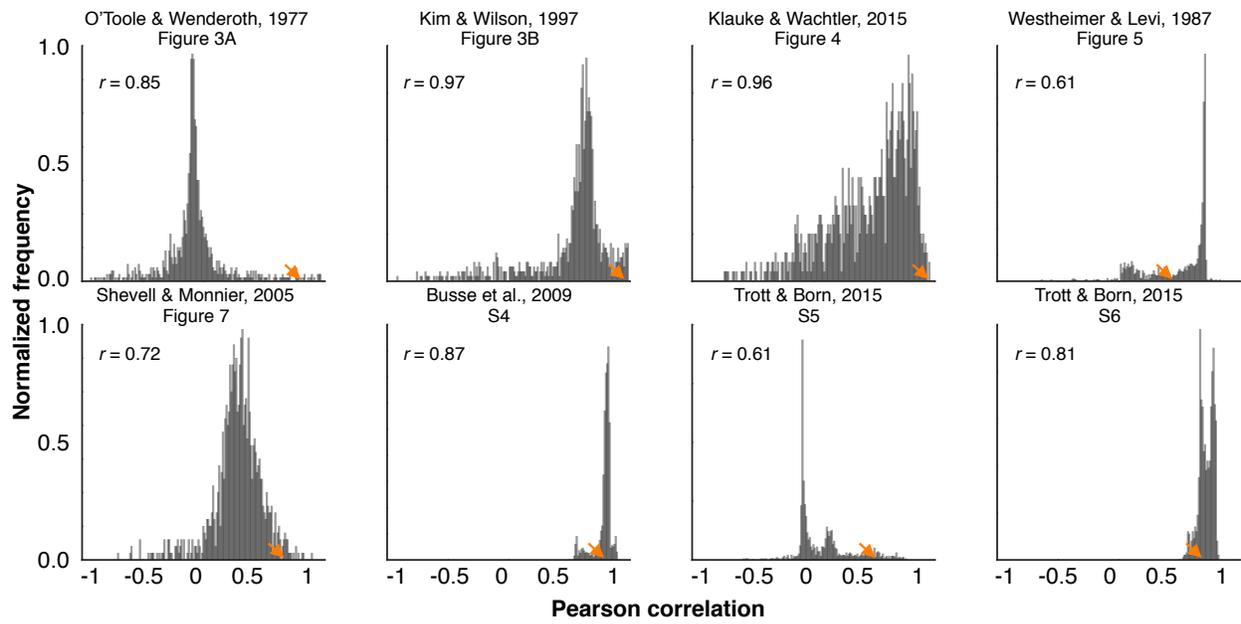


Figure S10. *Distributions of fit for the optimized full model when its strict one-to-one CRF-eCRF connectivity is relaxed.* Histograms depict the optimized full model’s ability to explain contextual phenomena (measured by Pearson correlation coefficients between model predictions and perceptual data) as the one-to-one connectivity of its CRF-eCRF connections are relaxed into weakly-tuned, “many-to-many” CRF-eCRF connectivity. For orientation, motion, and disparity phenomena, each eCRF unit had normally distributed connectivity, which was centered at a target unit tuning preference  $\theta_k$  with standard deviation  $\zeta$ . For the color domain we consider color-opponent model units with monotonic tuning curves and set  $w_{kk} = (\zeta\sqrt{2\pi})^{-1}$  and  $w_{jk} = \text{const.}$  (when  $j \neq k$ ; under the constraint that the weights sum up to  $r1$ ). The lesion screening procedure was used to evaluate the model’s ability to explain perception across a range of relaxed CRF-eCRF connectivity patterns. This involved randomizing  $\zeta$  over 1,000 iterations, and then measuring the model’s success in simulating the contextual phenomena with each pattern of connectivity. See 12 for details of the procedure for sampling  $\zeta$ . Pearson correlation values reported in the top left-hand corner of each panel denote performance of the optimized full model with strictly tuned connections on each problem. These values are also highlighted with orange arrows. Figure 6 is not included because it is computationally intractable for the lesion-screening.

952

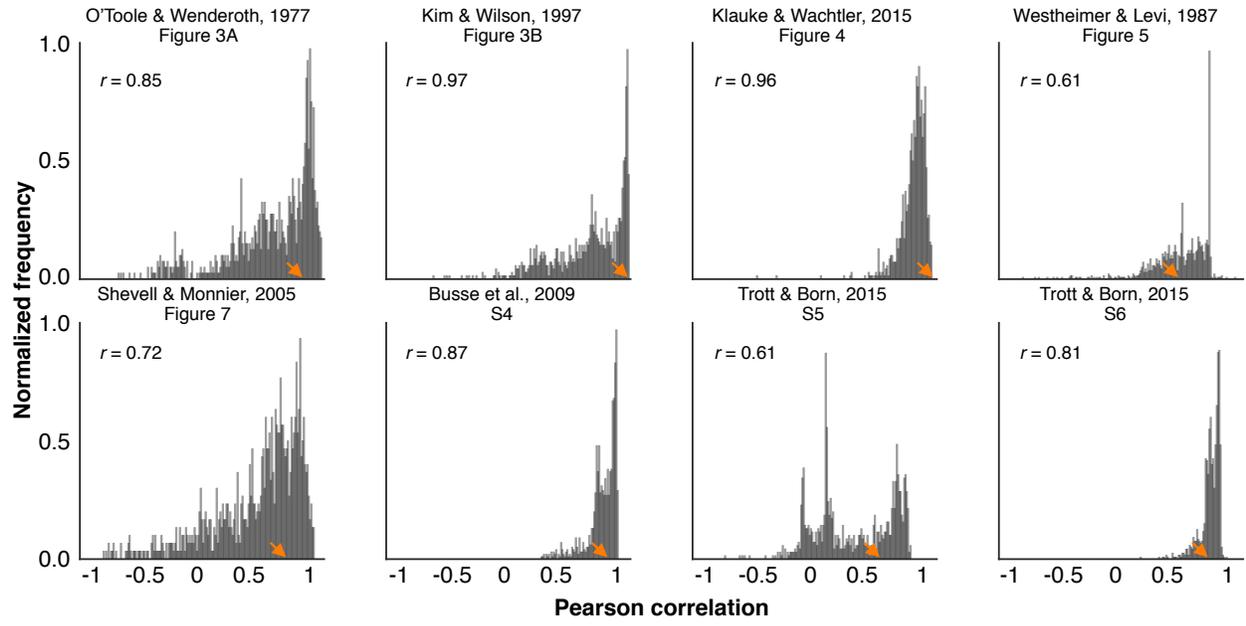


Figure S11. *Distributions of fit for the optimized full model derived by varying the tuning bandwidth of the input units' tuning curves.* Histograms depict the optimized full model's ability to explain perception (measured by Pearson correlations between model predictions and perceptual data) as the shapes of its tuning curves were randomized over 1,000 iterations for each contextual phenomenon. We explored a range of parameters around the hand-tuned values reported in the manuscript (see Eq. 12 for sampling details). For disparity (Cumming and Parker, 1997), motion direction (Albright et al., 1984) and orientation (Ringach et al., 1997) this involved resampling tuning curve bandwidth. For color opponent tuning this involved resampling the response threshold since we did not assume a bell-shape tuning curve for the modality. See 12 for details of the sampling procedure. Note that for the color domain the sampling range Pearson correlation values reported in the top left-hand corner of each panel denote performance of the optimized model with tuning bandwidth optimized for each experiment. These values are also highlighted with orange arrows. Figure 6 is not included because it is computationally intractable for the lesion-screening.

953

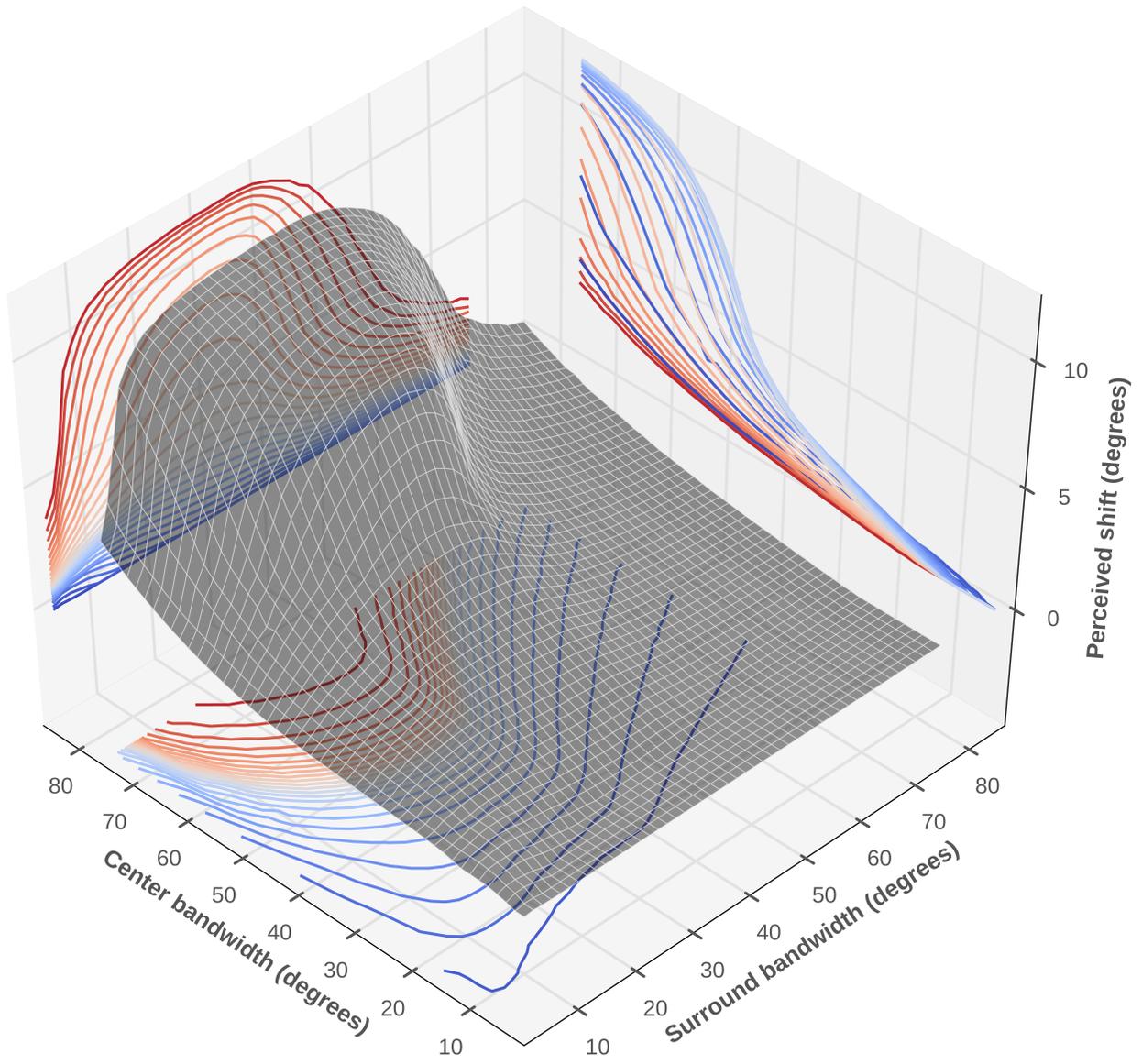


Figure S12. *Evolution of the repulsive shift as a function of tuning curve bandwidths.* We examined how the maximal amplitude of the repulsive orientation tilt effect (Figure 3A) changes as center and surround stimulus bandwidths are adjusted. The maximal amplitude is defined at the largest repulsive shift that is achieved as the orientation of the surround is systematically varied relative to the center.

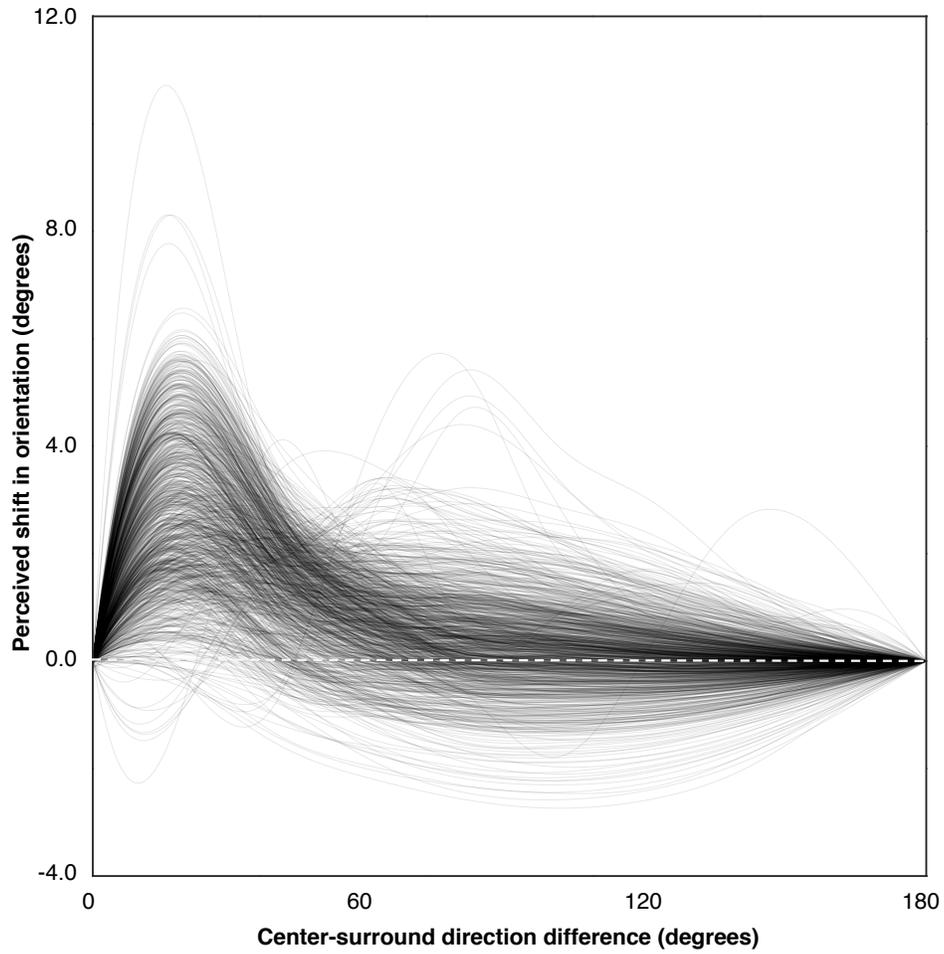


Figure S13. *Model instantiations derived from the lesion-screening procedure may account for some of the inter-participant variability observed experimentally.* Plot depicts full model behavior of the orientation tilt effect (Figure 3C) for every combination of parameters sampled (see Eq. 12 for sampling details).

955

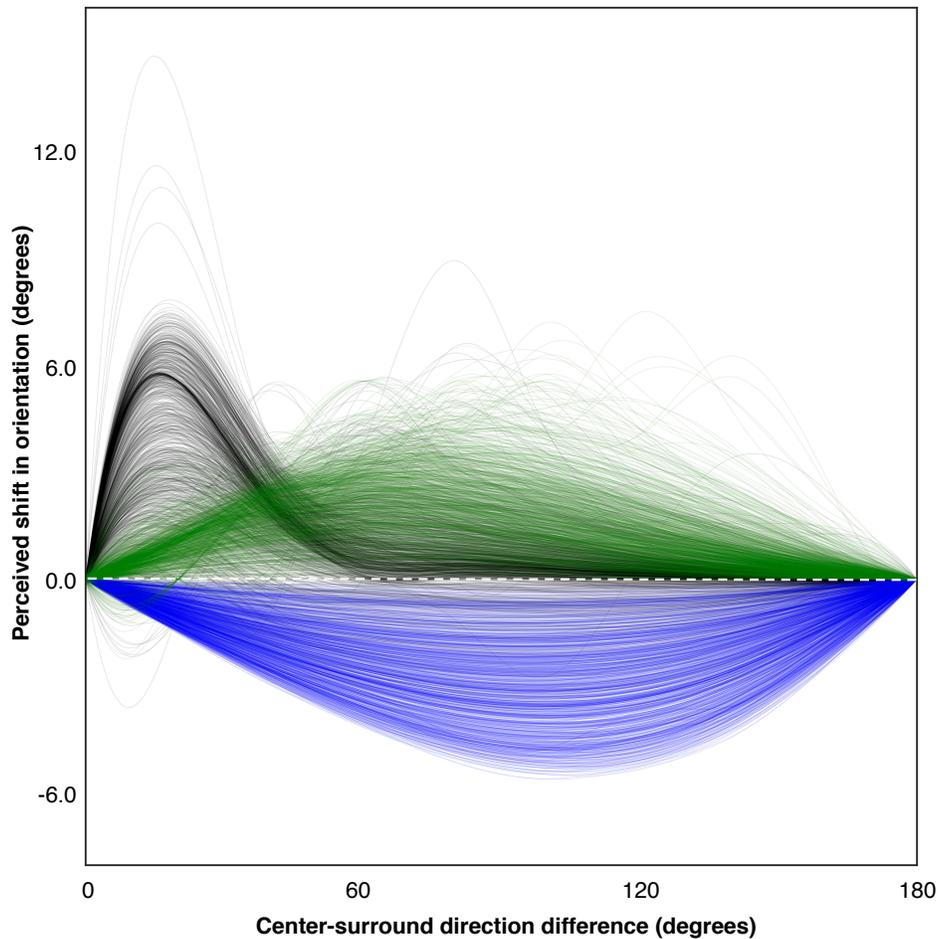


Figure S14. *Model instantiations with lesioned eCRF or divisive normalization mechanisms are unable to explain the orientation tilt effect.* Plot depicts simulations from models with lesions applied to either the facilitatory near eCRF (green), suppressive far eCRF (blue), or divisive normalization (black) on the orientation tilt effect (Figure 3C) for every sampled combination of parameters (see Eq. 12 for sampling details). Lesioning the facilitatory near eCRF destroys the attraction regime. Lesioning the suppressive far eCRF destroys the repulsion regime. Without  
 956 divisive normalization the model is unable to balance these two regimes.

## References

- 957  
958 Albright, T. D., Desimone, R., and Gross, C. G. (1984). Columnar organization of directionally  
959 selective cells in visual area MT of the macaque. *Journal of Neurophysiology*, 51(1):16–31.
- 960 Allman, J., Miezin, F., and McGuinness, E. (1985). Direction- and velocity-specific responses  
961 from beyond the classical receptive field in the middle temporal visual area (MT). *Perception*,  
962 14(2):105–26.
- 963 Angelucci, A., Levitt, J. B., and Lund, J. S. (2002a). Anatomical origins of the classical receptive  
964 field of single neurons in macaque visual cortical area V1. *Progress in Brain Research*,  
965 136:373–88.
- 966 Angelucci, A., Levitt, J. B., Walton, E. J. S., Hupe, J.-M., Bullier, J., and Lund, J. S. (2002b).  
967 Circuits for local and global signal integration in primary visual cortex. *Journal of Neuroscience*,  
968 22(19):8633–8646.
- 969 Angelucci, A. and Shushruth, S. (2013). Beyond the classical receptive field: surround modulation  
970 in primary visual cortex. In *The New Visual Neurosciences*, pages 425–444. MIT Press.
- 971 Anstis, S. (2006). White’s effect in lightness, color, and motion. In *Seeing Spatial Form*, pages  
972 91–100. Oxford University Press.
- 973 Anton-Erxleben, K., Stephan, V. M., and Treue, S. (2009). Attention reshapes center-surround  
974 receptive field structure in macaque cortical area MT. *Cereb. Cortex*, 19(10):2466–2478.
- 975 Baker, D. H. and Graf, E. W. (2010). Contextual effects in speed perception may occur at an early  
976 stage of processing. *Vision Research*, 50(2):193–201.

- 977 Bosking, W. H., Zhang, Y., Schofield, B., and Fitzpatrick, D. (1997). Orientation selectivity and the  
978 arrangement of horizontal connections in tree shrew striate cortex. *The Journal of Neuroscience*,  
979 17(6):2112–27.
- 980 Bradley, D. C. and Andersen, R. a. (1998). Center-surround antagonism based on disparity in  
981 primate area MT. *Journal of Neuroscience*, 18(18):7552–7565.
- 982 Briggs, F. and Usrey, W. M. (2011). Distinct mechanisms for size tuning in primate visual cortex.  
983 *The Journal of neuroscience*, 31(35):12644–12649.
- 984 Bringuier, V., Chavane, F., Glaeser L., and Frégnac, Y. (1999). Horizontal propagation of visual  
985 activity in the synaptic integration field of area 17 neurons. *Science*, 283(5402):695–699.
- 986 Busse, L., Wade, A. R., and Carandini, M. (2009). Representation of concurrent stimuli by  
987 population activity in visual cortex. *Neuron*, 64(6):931–42.
- 988 Carandini, M. (2012). From circuits to behavior: a bridge too far? *Nature neuroscience*,  
989 15(4):507–9.
- 990 Carandini, M. and Heeger, D. (2012). Normalization as a canonical neural computation. *Nature*  
991 *Reviews Neuroscience*, 13:51–62.
- 992 Carandini, M. and Heeger, D. J. (1994). Summation and division by neurons in primate visual  
993 cortex. *Science*, 264:1333–1336.
- 994 Cavanaugh, J. R., Bair, W., and Movshon, J. A. (2002). Selectivity and spatial distribution

995 of signals from the receptive field surround in macaque V1 neurons. *J Neurophysiol*,  
996 88(5):2547–2556.

997 Clifford, C. W. G. (2014). The tilt illusion: Phenomenology and functional implications. *Vision*  
998 *Research*, 104:3–11.

999 Coen-Cagli, R., Dayan, P., and Schwartz, O. (2012). Cortical surround interactions and perceptual  
1000 salience via natural scene statistics. *PLoS Computational Biology*, 8(3):e1002405.

1001 Conway, B. R. (2009). Color vision, cones, and color-coding in the cortex. *The Neuroscientist : a*  
1002 *review journal bringing neurobiology, neurology and psychiatry*, 15(3):274–90.

1003 Conway, B. R., Moeller, S., and Tsao, D. Y. (2007). Specialized color modules in macaque  
1004 extrastriate cortex. *Neuron*, 56(3):560–73.

1005 Cumming, B. G. and Parker, A. J. (1997). Responses of primary visual cortical neurons to  
1006 binocular disparity without depth perception. *Nature*, 389(6648):280–3.

1007 DeAngelis, G. C., Freeman, R. D., and Ohzawa, I. (1994). Length and width tuning of neurons in  
1008 the cat’s primary visual cortex. *Journal of neurophysiology*, 71(1):347–74.

1009 DeAngelis, G. C., Ghose, G. M., Ohzawa, I., and Freeman, R. D. (1999). Functional  
1010 Micro-Organization of Primary Visual Cortex: Receptive Field Analysis of Nearby Neurons.  
1011 *J. Neurosci.*, 19(10):4046–4064.

1012 DeAngelis, G. C. and Newsome, W. T. (1999). Organization of disparity-selective neurons in  
1013 macaque area MT. *J. Neurosci.*, 19(4):1398–1415.

- 1014 Douglas, R. J., Koch, C., Mahowald, M., Martin, K. A., and Suarez, H. H. (1995). Recurrent  
1015 excitation in neocortical circuits. *Science*, 269(5226):981–985.
- 1016 Dow, B. M. (2002). Orientation and color columns in monkey visual cortex. *Cerebral cortex*,  
1017 12(10):1005–1015.
- 1018 Eagleman, D. M. (2001). Visual illusions and neurobiology. *Nature Reviews Neuroscience*,  
1019 2(12):920–926.
- 1020 Georgopoulos, A. P., Schwartz, A. B., and Kettner, R. E. (1986). Neuronal population coding of  
1021 movement direction. *Science*, 233:1416–1419.
- 1022 Gibson, J. J. and Radner, M. (1937). Adaptation, after-effect and contrast in the perception of  
1023 tilted lines. *Journal of Experimental Psychology*, 20:453–467.
- 1024 Gilbert, C. D. and Li, W. (2013). Top-down influences on visual processing. *Nature reviews*.  
1025 *Neuroscience*, 14(5):350–63.
- 1026 Goddard, E., Clifford, C. W. G., and Solomon, S. G. (2008). Centre-surround effects on perceived  
1027 orientation in complex images. *Vision Research*, 48(12):1374–1382.
- 1028 Grossberg, S. and Todorović, D. (1988). Neural dynamics of 1-D and 2-D brightness perception: a  
1029 unified model of classical and recent phenomena. *Perception and psychophysics*, 43(3):241–277.
- 1030 Heeger, D. J. (1993). Modeling simple-cell direction selectivity with normalized, half-squared,  
1031 linear operators. *Journal of Neurophysiology*, 70(5):1885–1898.

- 1032 Hess, R. F., Hayes, A., and Field, D. J. (2003). Contour integration and cortical processing. *Journal*  
1033 *Of Physiology Paris*, 97(2-3):105–119.
- 1034 Hubel, D. and Wiesel, T. N. (1968). Receptive fields and functional architecture of monkey striate  
1035 cortex. *The Journal of physiology*, 195(1):215–43.
- 1036 Johnson, E. N., Hawken, M. J., and Shapley, R. (2001). The spatial transformation of color in the  
1037 primary visual cortex of the macaque monkey. *Nature neuroscience*, 4(4):409–16.
- 1038 Jones, H. E., Grieve, K. L., Wang, W., and Sillito, A. M. (2001). Surround suppression in primate  
1039 V1. *Journal of neurophysiology*, 86:2011–2028.
- 1040 Kim, J. and Wilson, H. R. (1997). Motion integration over space: interaction of the center and  
1041 surround motion. *Vision Research*, 37(8):991–1005.
- 1042 Klauke, S. and Wachtler, T. (2015). "Tilt" in color space: hue changes induced by chromatic  
1043 surrounds. *Journal of Vision*, 15(13):17.
- 1044 Kusunoki, M., Moutoussis, K., and Zeki, S. (2006). Effect of background colors on the tuning of  
1045 color-selective cells in monkey area V4. *Journal of neurophysiology*, 95(5):3047–59.
- 1046 Land, E. H. and McCann, J. J. (1971). Lightness and Retinex theory. *Journal of the Optical Society*  
1047 *of America*, 61(1):1–11.
- 1048 Lee, W.-c. A., Bonin, V., Reed, M., Graham, B. J., and Hood, G. (2016). Anatomy and function of  
1049 an excitatory network in the visual cortex. *Nature*, 532(7599):370–374.

1050 Levitt, J. B. and Lund, J. S. (1997). Contrast dependence of contextual effects in primate visual  
1051 cortex. *Nature*, 387(6628):73–76.

1052 Li, C. Y., Lei, J. J., and Yao, H. S. (1999). Shift in speed selectivity of visual cortical neurons: a  
1053 neural basis of perceived motion contrast. *Proceedings of the National Academy of Sciences of*  
1054 *the United States of America*, 96(7):4052–4056.

1055 Loomis, J. M. and Nakayama, K. (1973). A velocity analogue of brightness contrast. *Perception*,  
1056 2(4):425–428.

1057 Ma, W. J., Beck, J. M., Latham, P. E., and Pouget, A. (2006). Bayesian inference with probabilistic  
1058 population codes. *Nature Neuroscience*, 9(11):1432–8.

1059 Manassi, M., Lonchampt, S., Clarke, A., and Herzog, M. H. (2016). What crowding can tell us  
1060 about object representations. *J. Vis.*, 16(3):35.

1061 Marshak, W. and Sekuler, R. (1979). Mutual repulsion between moving visual targets. *Science*  
1062 *(New York, N.Y.)*, 205(4413):1399–401.

1063 Mollon, J. D. (2009). A neural basis for unique hues? *Current Biology*, 19(11):441–442.

1064 Monnier, P. and Shevell, S. K. (2003). Large shifts in color appearance from patterned chromatic  
1065 backgrounds. *Nature neuroscience*, 6(8):801–802.

1066 Murakami, I. and Shimojo, S. (1993). Motion capture changes to induced motion at higher  
1067 luminance contrasts, smaller eccentricities, and larger inducer sizes. *Vision Research*,  
1068 33(15):2091–2107.

- 1069 Murakami, I. and Shimojo, S. (1996). Assimilation-type and contrast-type bias of motion induced  
1070 by the surround in a random-dot display: evidence for center-surround antagonism. *Vision*  
1071 *Research*, 36(22):3629–3639.
- 1072 Nassi, J., Avery, M., Cetin, A., Roe, A., and Reynolds, J. (2015). Optogenetic Activation of  
1073 Normalization in Alert Macaque Visual Cortex. *Neuron*, 86(6):1504–1517.
- 1074 Norman, H. F., Norman, J. F., Todd, J. T., and Lindsey, D. T. (1996). Spatial interactions in  
1075 perceived speed. *Perception*, 25(7):815–830.
- 1076 O’Toole, B. and Wenderoth, P. (1977). The tilt illusion: repulsion and attraction effects in the  
1077 oblique meridian. *Vision research*, 17(3):367–374.
- 1078 Ozeki, H., Finn, I. M., Schaffer, E. S., Miller, K. D., and Ferster, D. (2009). Inhibitory stabilization  
1079 of the cortical network underlies visual surround suppression. *Neuron*, 62(4):578–92.
- 1080 Petrov, Y., Carandini, M., and McKee, S. (2005). Two distinct mechanisms of suppression in  
1081 human vision. *The Journal of neuroscience*, 25(38):8704–8707.
- 1082 Ringach, D. L., Hawken, M. J., and Shapley, R. (1997). Dynamics of orientation tuning in macaque  
1083 primary visual cortex. *Nature*, 387(6630):281–284.
- 1084 Ringach, D. L., Shapley, R. M., and Hawken, M. J. (2002). Orientation selectivity in macaque V1:  
1085 diversity and laminar dependence. *Journal of neuroscience*, 22(13):5639–5651.
- 1086 Rubin, D. B., Hooser, S. D. V., and Miller, K. D. (2015). The stabilized supralinear network : A  
1087 unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*, 85(1):1–51.

- 1088 Rust, N. C., Mante, V., Simoncelli, E. P., and Movshon, J. A. (2006). How MT cells analyze the  
1089 motion of visual patterns. *Nature Neuroscience*, 9(11):1421–31.
- 1090 Sceniak, M. P., Hawken, M. J., and Shapley, R. (2001). Visual spatial characterization of macaque  
1091 V1 neurons. *Journal of Neurophysiology*, 85(5):1873–1887.
- 1092 Sceniak, M. P., Ringach, D. L., Hawken, M. J., and Shapley, R. (1999). Contrast’s effect on spatial  
1093 summation by macaque V1 neurons. *Nature neuroscience*, 2(8):733–9.
- 1094 Schein, S. J. and Desimone, R. (1990). Spectral properties of V4 neurons in the macaque. *Journal*  
1095 *of Neuroscience*, 10(10):3369–3389.
- 1096 Schwabe, L., Obermayer, K., Angelucci, A., and Bress (2006). The role of feedback in shaping  
1097 the extra-classical receptive field of cortical neurons: A recurrent network model. *Journal of*  
1098 *Neuroscience*, 26(36):9117–9129.
- 1099 Schwartz, O., Hsu, A., and Dayan, P. (2007). Space and time in visual context. *Nat. Rev. Neurosci.*,  
1100 8:522–535.
- 1101 Sengpiel, F. (1997). Binocular rivalry: ambiguities resolved. *Current Biology*, 7:447–450.
- 1102 Sengpiel, F., Sen, A., and Blakemore, C. (1997). Characteristics of surround inhibition in cat area  
1103 17. *Experimental Brain Research*, 116(2):216–228.
- 1104 Series, P., Lorenceau, J., and Frégnac, Y. (2003). The “silent” surround of V1 receptive fields:  
1105 theory and experiments. *Journal of physiology-Paris*, 97:453–474.

- 1106 Setić, M. and Domijan, D. (2008). Modeling the top-down influences on the lateral interactions in  
1107 the visual cortex. *Brain Res.*, 1225:86–101.
- 1108 Shevell, S. K. and Monnier, P. (2005). Color shifts from S-cone patterned backgrounds: Contrast  
1109 sensitivity and spatial frequency selectivity. *Vision Research*, 45:1147–1154.
- 1110 Shushruth, S. and Ichida, J. (2009). Comparison of spatial summation properties of neurons in  
1111 macaque V1 and V2. *Journal of neurophysiology*, 1:2069–2083.
- 1112 Shushruth, S., Mangapathy, P., Ichida, J. M., Bressloff, P. C., Schwabe, L., and Angelucci, A.  
1113 (2012). Strong recurrent networks compute the orientation tuning of surround modulation in the  
1114 primate primary visual cortex. *Journal of Neuroscience*, 32(1):308–321.
- 1115 Sillito, A. M., Grieve, K. L., Jones, H. E., Cudeiro, J., and Davis, J. (1995). Visual cortical  
1116 mechanisms detecting focal orientation discontinuities. *Nature*, 378:492–496.
- 1117 Sincich, L. and Horton, J. (2005). The circuitry of V1 and V2: integration of color, form, and  
1118 motion. *Annu. Rev. Neurosci.*, 28(1):303–326.
- 1119 Sit, Y. F., Chen, Y., Geisler, W. S., Miikkulainen, R., and Seidemann, E. (2009). Complex  
1120 Dynamics of V1 Population Responses Explained by a Simple Gain-Control Model. *Neuron*,  
1121 64(6):943–956.
- 1122 Smith, V. C., Jin, P. Q., and Pokorny, J. (2001). The role of spatial frequency in color induction.  
1123 *Vision Research*, 41(8):1007–1021.

- 1124 Stemmler, M., Usher, M., and Niebur, E. (1995). Lateral interactions in primary visual cortex: a  
1125 model bridging physiology and psychophysics. *Science*, 269(5232):1877–1880.
- 1126 Stepanyants, A., Hirsch, J. A., Martinez, L. M., Kisvarday, Z. F., Ferecsko, A. S., and Chklovskii,  
1127 D. B. (2008). Local potential connectivity in cat primary visual cortex. *Cerebral Cortex*,  
1128 18(1):13–28.
- 1129 Tanaka, H. and Ohzawa, I. (2009). Surround suppression of V1 neurons mediates orientation-based  
1130 representation of high-order visual features. *Journal of neurophysiology*, 101(3):1444–62.
- 1131 Thomas, O. M., Cumming, B. G., and Parker, a. J. (2002). A specialization for relative disparity  
1132 in V2. *Nature neuroscience*, 5(5):472–8.
- 1133 Trott, A. R. and Born, R. T. (2015). Input-gain control produces feature-specific surround  
1134 suppression. *Journal of neuroscience*, 35(12):4973–1982.
- 1135 Tsotsos, J. K., Culhane, S., and Cutzu, F. (2001). Visual attention and cortical circuits. In Braun,  
1136 J., Koch, C., and Davis, J. L., editors, *From Theoretical Foundations to a Hierarchical Circuit*  
1137 *for Selective Attention*, pages 285–306. The MIT presss.
- 1138 Vinck, M., Womelsdorf, T., and Fries, P. (2013). Gamma-Band synchronization and information  
1139 transmission. In *Principles of Neural Coding*, pages 449–470. CRC Press.
- 1140 Wachtler, T., Sejnowski, T. J., and Albright, T. D. (2003). Representation of color stimuli in awake  
1141 macaque primary visual cortex. *Neuron*, 37(4):681–91.
- 1142 Weliky, M., Kandler, K., Fitzpatrick, D., and Katz, L. C. (1995). Patterns of excitation and

- 1143 inhibition evoked by horizontal connections in visual cortex share a common relationship to  
1144 orientation columns. *Neuron*, 15(3):541–552.
- 1145 Westheimer, G. (1986). Spatial interaction in the domain of disparity signals in human stereoscopic  
1146 vision. *The Journal of physiology*, 370:619–29.
- 1147 Westheimer, G. and Levi, D. M. (1987). Depth attraction and repulsion of disparate foveal stimuli.  
1148 *Vision Research*, 27(8):1361–1368.
- 1149 White, M. (1979). A new effect of pattern on perceived lightness. *Perception*, 8(4):413–416.
- 1150 White, M. (1981). The effect of the nature of the surround on the perceived lightness of grey bars  
1151 within square-wave test gratings. *Perception*, 10(2):215–230.
- 1152 Zhang, J., Barhomi, Y., and Serre, T. (2012). A new biologically inspired color image descriptor.  
1153 In *European Conference on Computer Vision*, pages 312–324.
- 1154 Zhu, M. and Rozell, C. J. (2013). Visual nonclassical receptive field effects emerge from Sparse  
1155 coding in a dynamical system. *PLoS Computational Biology*, 9(8):1–15.