# A simple circuit model of visual cortex explains neural and behavioral aspects of attention

Grace W. Lindsay[a,*], Daniel B. Rubin[b,*], Kenneth D. Miller[c]

[a]*Gatsby Computational Neuroscience Unit, Sainsbury Wellcome Centre, University College London, London, UK*
[b]*Department of Neurology, Massachusetts General Hospital, Harvard Medical School, Boston MA*
[c]*Center for Theoretical Neuroscience, College of Physicians and Surgeons, Mortimer B. Zuckerman Mind Brain Behaviour Institute, Swartz Program in Theoretical Neuroscience, Kavli Institute for Brain Science, New York, Department of Neuroscience, Columbia University, New York, United States*

## Abstract

Selective visual attention modulates neural activity in the visual system and leads to enhanced performance on difficult visual tasks. Here, we use an existing circuit model of visual cortex, known as the stabilized supralinear network, to demonstrate that many neural correlates of attention can arise from simple circuit mechanisms. Using different variants of the model we replicate results from studies of both feature and spatial attention. In addition to firing rate changes, we also replicate findings regarding how attention impacts trial-to-trial variability. Finally, we expand this circuit model into an architecture that can perform visual tasks in order to show that these neural effects can enhance detection performance. This work advances our understanding of the physical underpinnings of attention.

*Keywords:* Attention, Normalization, Neural Networks

## 1. Introduction

²    When an animal knows in advance what features or locations in the visual ³ scene will be relevant for completing its goals, selective top-down attention

---

*These authors contributed equally to the work. Corresponding author: Grace W. Lindsay, gracewlindsay@gmail.com

can be deployed. This attention has been shown to have a powerful modulatory effect on both task performance and neuronal responses, and changes in the latter can often be powerful predictors of the former (Ress et al., 2000).

Numerous specific impacts of attention on neural activity have been identified, including changes in firing rates, trial-to-trial variability, and noise correlations (Treue and Maunsell, 1999; Treue and Martinez Trujillo, 1999; Cohen and Maunsell, 2009). Looking at the impact of attention on tuning curves, attention to a preferred stimulus is known to scale up the responses to all stimuli; conversely, attention to a non-preferred stimulus scales responses down (Martinez-Trujillo and Treue, 2004). This enhancement has been shown to be a largely multiplicative increase in neuronal gain (Treue and Martinez Trujillo, 1999). A similar percentage change occurs in the firing rates of excitatory and inhibitory neurons (Mitchell et al., 2007).

Many of attention's impacts on firing rates can be understood in the context of the normalization model of attention (Reynolds and Heeger, 2009; Lee and Maunsell, 2009; Ghose, 2009; Boynton, 2009). This model builds off the canonical computation of normalization observed in multiple places in the visual system as well as other brain areas (Carandini and Heeger, 2012). In the absence of attention, a neuron's firing rate can be predicted by a divisive normalization equation: stimuli with the preferred features and in the classical receptive field of the neuron form the numerator (known as the "stimulus drive"), and the denominator is a function of a less-selective suppressive drive that includes surround locations and non-preferred features as well. Under the normalization model of attention, attention provides a biasing effect that amplifies the drive coming from the attended stimulus.

This model captures how attention can, when two stimuli are present, shift responses to those of the attended stimulus alone. For example, when a preferred and non-preferred stimulus are both presented to the receptive field of a V4 neuron, the cell's response is intermediate between the responses evoked by each stimulus alone. By attending to either the preferred or non-preferred stimulus, the response is shifted towards the response evoked by the attended stimulus alone (Reynolds and Desimone, 2003). Similarly, attention to a stimulus in the suppressive surround of a V4 neuron increases the suppression induced, whereas attention to the center reduces the suppression (Sundberg et al., 2009). The normalization model of attention also captures how attention increases contrast gain or response gain, respectively, depending on whether the attention is over a larger or smaller cortical area than the stimulus input (Reynolds and Heeger, 2009).
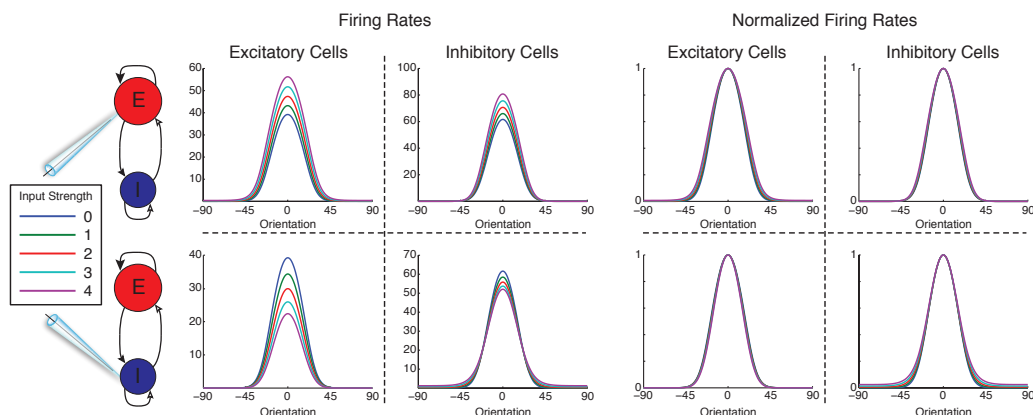
2

Figure 1:

**Expansive nonlinearity and balanced amplification yield multiplicative scaling.** We consider a simple two-unit nonlinear SSN model, with one excitatory (E) cell and one inhibitory (I) cell (Methods 4.1.4). We drove both cells with a series of feedforward inputs, whose strengths varied as a function of "orientation" to generate "tuning curves". While driving the cells with this feedforward input, an additional constant input of one of four varying strengths (indicated by color legend at left) was added to either the E or the I cell. With increasing input to the E cell, both E and I rates are scaled up, whereas with increasing input to the I cells, both E and I rates are scaled down. Normalizing each curve by its maximum reveals that the gain change is almost exclusively multiplicative.

Beyond changes in firing rates described by the normalization model of attention, attention also decreases trial-to-trial variability and noise correlations across neuron pairs (Cohen and Maunsell, 2009; Mitchell et al., 2007). We have previously shown that a simple model of cortical circuitry— known as the stabilized supralinear network (SSN) (Ahmadian et al., 2013)— can account for a wide set of phenomena described by the normalization model, including feature normalization and surround suppression and their nonlinear dependencies on contrast (Rubin et al., 2015). It also accounts for the suppression of correlated variability by a stimulus (Hennequin et al., 2018). The network assumes expansive or supralinear input/output functions for the individual units. As described in (Ahmadian et al., 2013; Rubin et al., 2015; Ahmadian and Miller, 2019), this causes effective synaptic strengths between units, which are proportional to the postsynaptic neuron's gain –

3

its change in firing rate for a given change in input – to grow with increasing postsynaptic activation. The growth of excitatory-to-excitatory effective connections leads to potential instability, but with sufficiently strong feedback inhibition the network remains stable. However, this stabilization occurs through the network dynamically "loosely balancing" its inputs, so that the recurrent input largely cancels the feedforward input, leaving a residual net input that grows sublinearly as a function of the feedforward input. (The balancing is "loose" because the residual input after cancellation is comparable in size to the factors that cancel, Ahmadian and Miller, 2019.) This cancellation of feedforward input through increasingly strong inhibitory stabilization leads to the normalization and variability suppression effects just described.

The SSN has strong recurrent excitation stabilized by strong feedback inhibition and exhibits "balanced amplification" (Hennequin et al., 2018; Murphy and Miller, 2009): small inputs biased toward either excitatory (or inhibitory) cells drive large increases (or decreases) in both excitatory and inhibitory firing rates. We hypothesized that attentional modulation acts through the same balanced amplification and recurrent "loose balancing" mechanisms that implement feature normalization and surround suppression. Here we show that this model can indeed account for many of the neural effects of attention observed in visual cortex.

Finally, in addition to replicating neural effects, we also use this model to show how changes in neural activity can enhance performance. Previous work (Lindsay and Miller, 2018) used a deep convolutional neural network (CNN) as a model of the visual system to show how neural changes associated with attention enhance performance on a challenging visual detection task. Here, we put our circuit model into a convolutional architecture to create a model that connects low-level circuitry with behavioral outputs. This model (dubbed the SSN-CNN) replicates both the neural impacts of attention as well as the performance enhancements.

## 2. Results

We employ four instantiations of our model of visual cortex to replicate the neural effects of attention. The details of all of these models have been described previously, and are included in the Methods section. All four models feature strongly recurrently connected excitatory and inhibitory neurons with a supralinear neuronal input-output nonlinearity. The four models differ

4

only in the dimension of stimulus space over which the neurons are arranged and the spatial arrangement and strengths of the connections between neurons. In the simplest model, we consider a single pair of excitatory and inhibitory neurons (Figure 1). The two slightly more complex models represent populations of neurons either arranged around a ring, with position on the ring interpreted as preferred orientation of cells with a similar retinotopic receptive field (RF) position(Methods 4.1.1, Figure 2), or on a line, with position on the line interpreted as retinotopic RF position of cells with similar preferred orientation (Methods 4.1.2, Figure 15). The most complex model has a 2-dimensional representation of retinotopic space on which is superimposed a map of preferred orientations. In this model, neurons make connections as probabilistic functions of difference in stimulus preference over the three dimensions of stimulus quality: two spatial dimensions and orientation (Methods 4.1.3).

We note that the suppression of response to a preferred orientation by simultaneous presentation of an orthogonal orientation or "mask" ("cross-orientation suppression") in V1 is largely mediated by nonlinear changes in the pattern of thalamic firing induced by the mask, rather than by nonlinear V1 integration (Priebe and Ferster, 2006; Li et al., 2006), although there is a component mediated by V1 as shown by suppression arising when the two stimuli are presented to different eyes (Sengpiel and Vorobyov, 2005). In our models, the inputs to the model cortex are assumed to sum linearly, so that all nonlinear behavior arises from cortical processing. We typically refer to different competing stimuli presented within an RF as "orientations", but this should be understood to model cortical processing given linear summation of inputs induced by two stimuli, rather than the literal phenomenon of V1 cross-orientation suppression.

In all instantiations, attention is modeled as a small additional excitatory input biased towards the excitatory cells within the specified locus of attention. As a secondary test, we also re-ran all simulations with attention instead modeled as a small inhibitory input towards the inhibitory cells (resulting in a disinhibition of locally-connected excitatory cells). Results were qualitatively similar, with a few notable exceptions discussed below.

To investigate the impact of neural activity changes on performance, we also incorporated one of these circuit models—the ring model—into a convolutional neural network architecture (Methods 4.3). This allowed us to demonstrate that the application of attention to our circuit model can increase performance on a challenging visual detection task.
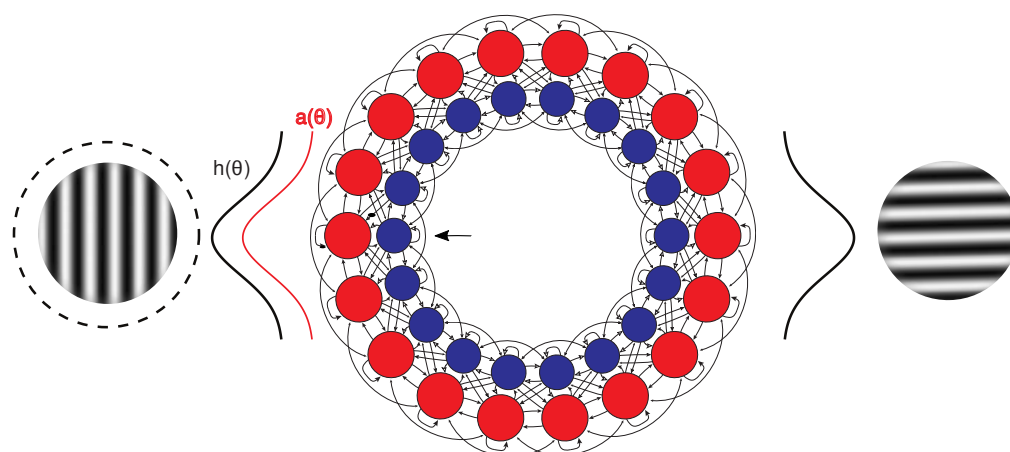
Figure 2:

**A ring model of attention.** The ring model represents different features (*e.g.*, preferred orientation) at a single location in visual space. At each location on the ring, a pair of excitatory (red) and inhibitory (blue) cells exist. Oriented stimuli are modeled as Gaussians centered at a particular location on the ring (black curves). Attention to one of the stimuli (indicated by dashed circle around it) is modeled as an additional Gaussian input biased towards the excitatory subpopulation at the center of the locus of attention (red curve). In this example, recording from the E-I pair indicated with the arrow would correspond to the cyan line in Figure 3.

### 2.1. Basic mechanism of the model

The balanced amplification model (Murphy and Miller, 2009) demonstrates that in a network with strong recurrent connectivity, small changes in the difference between E and I activity can drive large changes in the sum of the activity. Previously, we have used this mechanism to produce models of contextual modulation that capture the experimental observation that, during surround suppression, both E and I firing rates are suppressed (Ozeki et al., 2009). Within a locus of attention, however, the opposite effect is observed: both E and I firing rates are enhanced (Mitchell et al., 2007).

In a network wherein neurons are described by a supralinear nonlinearity, a bias in the input towards E or I shifts the responses of both cells up or down (respectively), and the resulting change can be almost exclusively multiplicative (Figure 1). Thus we hypothesize that this simple, intrinsic
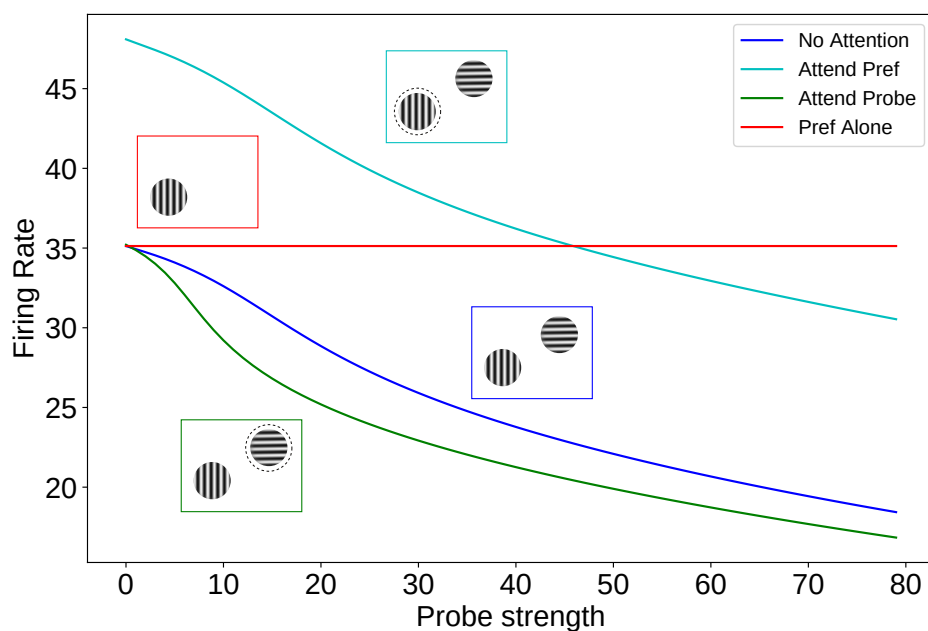
6

Figure 3: **Attention enhances the suppressive effect of non-preferred stimuli** A stimulus of preferred orientation was shown to a cell in the ring model. An orthogonally oriented stimulus was presented along with the preferred stimulus, and the strength of the non-preferred "probe" was varied (blue line). The test was then repeated with attention (indicated by dashed circle around stimulus) directed towards either the preferred stimulus (cyan) or the probe stimulus (green). When attention was directed towards the preferred stimulus, suppression was decreased. When attention was directed to the probe stimulus, suppression was enhanced.

142 form of amplification may be sufficient to account for the observed effects
143 of attention on visual cortical circuits. We now incorporate this simple E-I
144 pair into a broader recurrent circuit and consider several recent experimental
145 results on attention in visual cortex.

146 *2.2. Attention influences stimulus interactions*

147 *2.2.1. Impact of feature attention*

148    In several regions of visual cortex, attention to one of multiple stimuli
149 presented within the receptive field of a neuron can shift the response of that

7

neuron towards the response evoked by the attended stimulus alone. This was shown by Reynolds and Desimone (2003), who probed the responses of V4 neurons with preferred and non-preferred stimuli, presented either alone or together in the receptive field of a single neuron. They found that in the simultaneous presentation condition, attending to a non-preferred stimulus caused a relative suppression compared to an attend-away condition, whereas attending to the preferred stimulus boosted the response. To simulate this experiment, we recorded the response of a cell to a strong stimulus of preferred orientation in the ring model (for details of attention experiments see Methods 4.2). We then added a non-preferred stimulus at the orthogonal orientation to the ring (schematized in Figure 2) and systematically varied the strength of this "probe" stimulus. As expected, the addition of the non-preferred probe was always suppressive, and with increasing probe strength suppression was increased (Figure 3, blue line). We then repeated the same test with attention directed either towards the preferred stimulus (cyan) or the probe stimulus (green). When attention was directed towards the preferred stimulus, the amount of suppression was decreased. When attention was directed to the probe stimulus, suppression was enhanced.

In a related experiment, Treue and Martinez-Trujillo (1999) recorded from a neuron in area MT while presenting two stimuli to the neuron's receptive field. One of the stimuli was always moving in a non-preferred direction, while the direction of the other stimulus was systematically varied. Compared to an attend-away condition, responses of MT neurons were relatively suppressed at all stimulus directions when attention was directed towards the non-preferred stimulus, but relatively enhanced when attending towards the varying stimulus. We find the same result if we repeat this test in our ring model (Figure 4). Like Treue and Martinez-Trujillo (1999), the change we observe occurred without a substantial change in the width of tuning, indicating a mainly multiplicative scaling (Figure 4, inset).

Note that in Figures 3 and 4 the same strength of attention is applied in all circumstances, however attention applied to a non-preferred stimulus has a weaker impact on firing rates. In our model, attention applied to a cell's preferred stimulus means additional excitatory inputs to the cell in question. Attention to an orthogonal stimulus only impacts the recorded cell indirectly through recurrent connections, leading to a weaker effect. Experimentally, the magnitude of firing rate changes has been found to be weaker when attention is applied to a non-preferred stimulus compared to a preferred one (Treue and Maunsell, 1999).
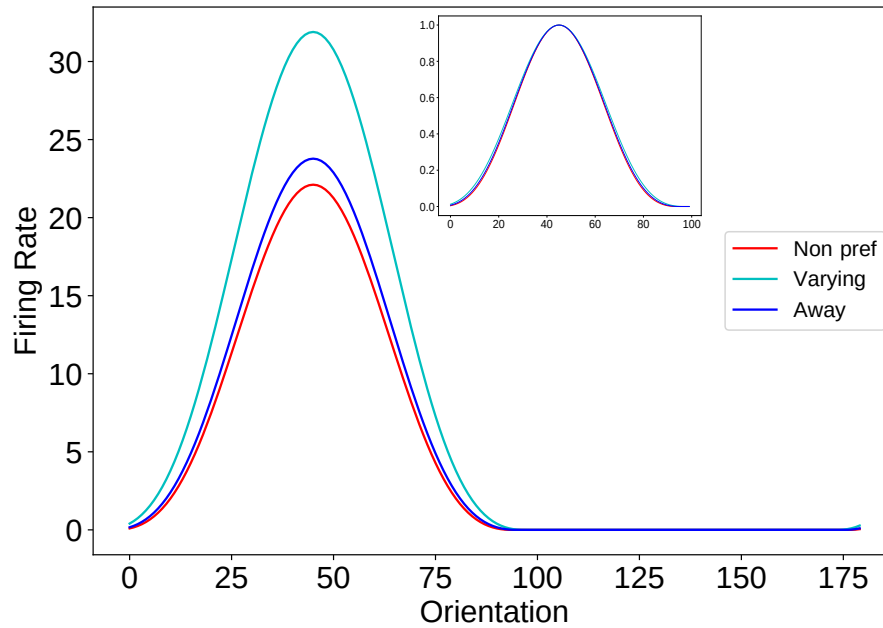
8

Figure 4:

**Attention scales tuning multiplicatively.** In the presence of a non-preferred probe stimulus, we varied the orientation of a test stimulus between 0° and 180°, while recording from the cell at 45° and attending either to the non-preferred probe (red), the varying stimulus (cyan), or away (blue). Attention produced an almost exclusively multiplicative change in response. Normalized responses are shown in the inset. There was virtually no change in tuning width, as observed experimentally (Treue and Martinez Trujillo, 1999).

9

### 2.2.2. Correlation between feature attention and normalization

Several groups have considered the mechanistic relationship between attention and cortical normalization (Reynolds and Heeger, 2009; Lee and Maunsell, 2009; Ni et al., 2012). In a recent study exploring the variability in the strength of attentional modulation, Ni and collegues demonstrated that neurons vary in the degree to which their responses are normalized by the presence of an orthogonal, non-preferred stimulus in the receptive field. They further show that the degree of normalization a cell demonstrates (or in their terminology, the broadness of the "tuning" of normalization – quantified by a normalization modulation index) is highly correlated with the extent to which attention modulates the response to the cell. To simulate this experiment, we employed our 2-D model of visual cortex designed to reproduce both the mean effects as well as a realistic degree of variability in responses. In this simulation, excitatory cells were selected at random from the population. For each cell, a high contrast stimulus of preferred orientation was presented. An orthogonal stimulus of the same size, position, and strength (the "null" stimulus) was then presented, and then the preferred and orthogonal stimuli were presented together. The firing rate response in each of the three stimulus conditions was recorded, and the Normalization Modulation Index was calculated as: NMI = [(r(Preferred) - r(Null)) - (r(Both - r(Null))]/[(r(Preferred) - r(Null)) + (r(Both - r(Null))]. An NMI of 0.33 corresponds to averaging of the two stimuli, whereas an NMI of 0 is considered a "winner take all" response (the response to the pair is the same as the response to the preferred stimulus alone). In the terminology of Ni et al., cells with highly tuned normalization have an NMI closer to 0 (Ni et al., 2012). The paired presentations were then repeated (showing both preferred + null together) with attention directed towards either the preferred or null stimulus. Attention was applied to the E cells in the position, size, and orientation of either the preferred or null stimulus. An Attentional Modulation Index was then calculated as: AMI = (r(Attend Preferred) - r(Attend Null))/(r(Attend Preferred) + r(Attend Null)). As was observed experimentally, there is a wide range of NMIs and AMIs, and the NMI and AMI of cells are highly correlated (Figure 5).

### 2.2.3. Impact of spatial attention

The previously discussed experiments studied the response of neurons to pairs of stimuli presented within the same receptive field. However, attention has also been shown to modulate the effect of stimuli presented in
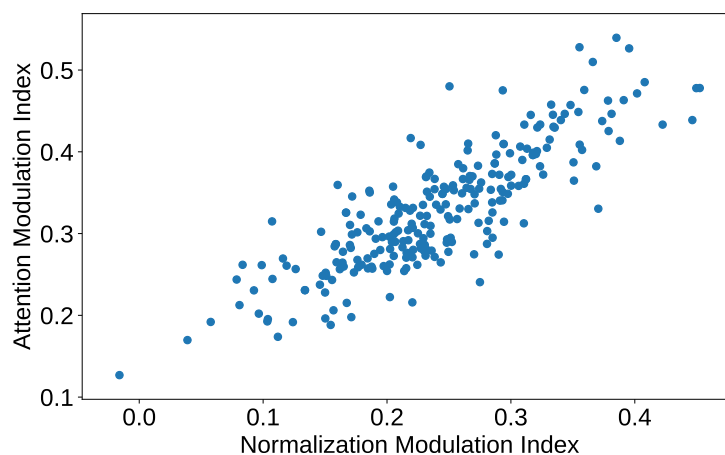
10

Figure 5:

**Normalization strength and attentional modulation are positively correlated.** Normalization Modulation Indices are plotted against the Attention Modulation Indices for all 250 cells sampled from the 2-D model. Correlation coefficient: 0.84. See text for details.

the receptive field surround. Sundberg et al. (2009) found that in V4, the strength of surround suppression could be either increased or decreased by attending specifically to the surround or center stimulus. To simulate this experiment, we next employed our line model used to simulate spatial contextual interactions. Pairs of E and I cells are arranged along a one-dimensional lattice representing an axis of retinotopic space, with recurrent excitatory connections that decrease as a function of retinotopic/cortical distance. It is assumed that the cells share preferred features. A stimulus was presented to the cell in the center of the lattice, in the presence of a suppressive surround stimulus. Attention was then directed to either the center or surround stimulus. Attention to the center decreased the strength of surround suppression (pushing firing rates towards those when the stimulus is presented alone), while attention to the surround enhanced surround suppression (Figure 6).

We simulated this experiment in the 2-D model as well. 100 neurons were randomly selected from the network. For each neuron, we measured the response to a strong stimulus of preferred orientation centered on the receptive field, and then added a strong stimulus of the same orientation to the surround. The response to the cell was measured in the absence
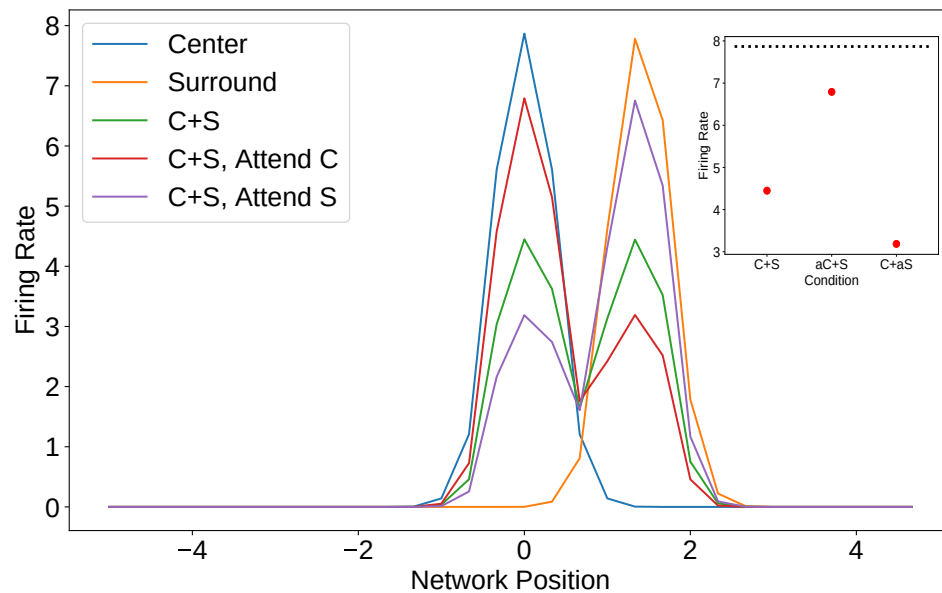
11

Figure 6:

**Attention modulates the strength of surround suppression.** A stimulus was shown in the receptive field of the neuron at position 0. A stimulus of equal strength and size was then placed in the surround, and the response was recorded from neurons in the vicinity. Attention was then directed either to the center or surround stimulus. In the main figure, the E cell activity across the network is shown in response to the center stimulus alone, the surround stimulus alone, the center and surround stimuli shown together, the center and surround stimuli with attention directed towards the center, and the center and surround stimuli with attention directed towards the surround. The inset demonstrates the activity at the center E cell – the dashed line is the response to the center stimulus alone, and the three dots show the response to the center and surround presented together, either with no attention, with attention directed towards the center, or with attention directed towards the surround.

of an attentional input (the "Attend Away" condition), as well as with an attentional input directed towards the center or surround stimulus. As was observed experimentally, attending to the surround boosted the amount of surround suppression, whereas attending to the center greatly weakened the
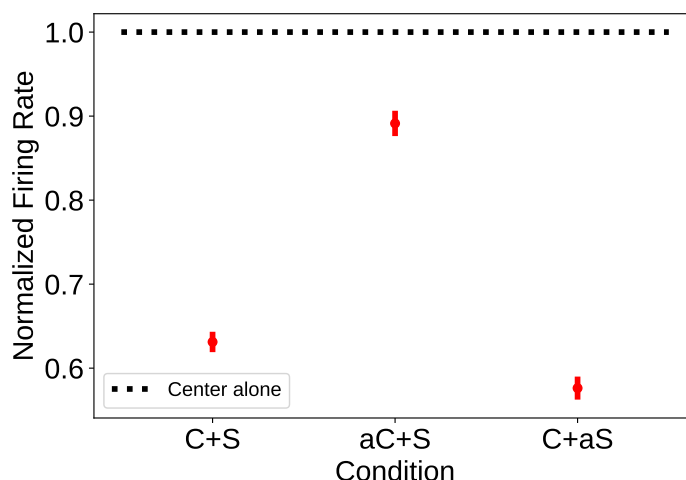
12

Figure 7:

**Attention modulates the strength of surround suppression in the large scale model.** A stimulus of preferred orientation was shown to a randomly selected cell. A stimulus with the same orientation and strength was placed in the surround, and the response was recorded. Attention was then directed either to the center or surround stimulus. The mean responses relative to the center alone is shown for a sample of 100 neurons from the 2-D model. Error bars indicate the standard error of the mean. All three response groups are significantly different from each other at $p < .005$ (student's t-test).

surround suppression (Figure 7, compare the results of the 2-D model to the inset of Figure 6).

## 2.3. Experimental paradigm alters the impact of attention

### 2.3.1. Effect on contrast and response gain

All of the experiments and simulations discussed thus far demonstrate that attention produces a gain change in the firing rate of neurons within the locus of attention. The quality of this gain change, however, can be strongly influenced by the relative sizes of the stimulus and the attentional field. Reynolds and Heeger (2009) (their Figure 3) found in their normalization model of attention that when attention is directed to a relatively large area, the effect on the response to a small stimulus should be predominantly a change in "contrast-gain", such that cells respond to stimuli as if they were effectively at higher contrast. This would be seen as a leftward shift in a

13

contrast-response curve for a stimulus, with relatively little change in the maximum firing rate. For a large stimulus and a small attentional field, they instead predict a change in "response-gain", such that all responses are scaled multiplicatively.

Here we again employ the one-dimensional spatial line network model to study the two different effects of attention described by Reynolds and Heeger (2009). Attention was still modeled as a small additional input only to excitatory cells over a defined spatial area, and we calculated "contrast response curves" with and without attention. (Note that what we call "contrast" is actually external input strength, *i.e.* the parameter $c$ in Eq. 3; in reality, external input strength, as measured by thalamic input firing rate, is a monotonic but nonlinear, saturating function of stimulus contrast, (*e.g.* Sclar, 1987; Sclar et al., 1990).) To quantify changes in the contrast response properties, we fit each curve to a standard Naka-Rushton equation (Naka and Rushton, 1966):

$$R(c) = R_{max} \left( \frac{c^n}{c_{50}^n + c^n} \right) \tag{1}$$

where $R_{max}$ is the plateau firing rate, $n$ describes the steepness of the contrast response curve, and $c_{50}$ is the strength of the stimulus at which the response is 50% of its maximum. In our fitting procedure, the value of $n$ is discovered for the no-attention condition, and held at that value when fitting the attended condition.

With a large attentional field and small stimulus, the effect of attention was predominantly a leftward shift in the contrast-response function, as predicted by the model of Reynolds and Heeger (2009). We quantified this change in "contrast gain" as the difference in the $c_{50}$ parameters of the contrast response curves produced with and without attention (Figure 8A). We compared this to the "response gain", which we quantify as the ratio of $R_{max}$ parameters with and without attention. With a large stimulus and small attentional field, the effect of attention was reversed: there was little change in the contrast gain, and a much larger change in the response gain (Figure 8B). The dashed lines in either figure show the percent change in firing rate induced by attention. With a change in contrast gain there is little change in firing at the largest contrast, but this is not true for a change in response gain.

While Reynolds and Heeger (2009) showed this property in their descriptive model of attention, conditions that produce changes in contrast or re-
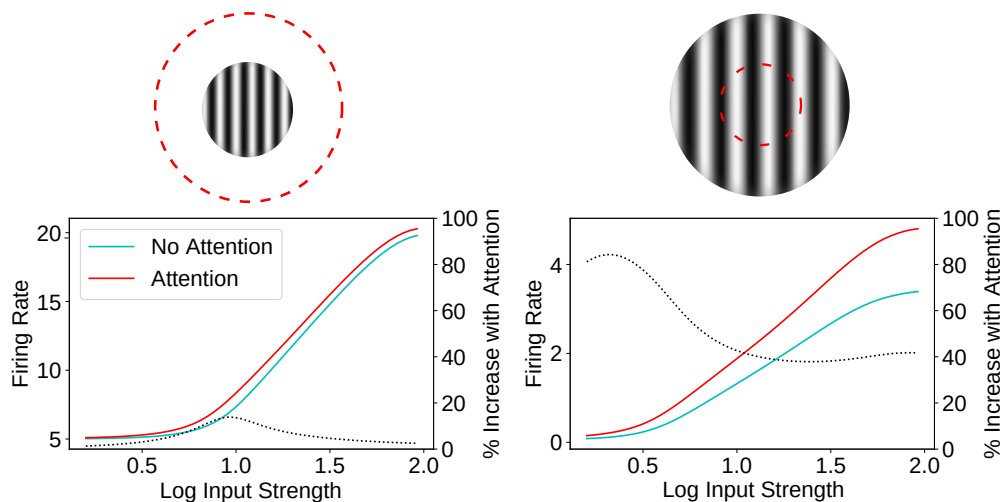
14

Figure 8:

**The qualitative effect of attention depends on the relative sizes of the attentional and stimulus fields.** Here we used the spatial line model to study the two different effects of attention, as described by Reynolds and Heeger (2009), Figure 3. Contrast response curves were calculated by varying the input strength logarithmically (base 10) in the presence (red curves) and absence (cyan curves) of attention. Left: with a large attentional field (red dashed circle) and small stimulus, the impact of attention was largely on contrast gain, defined as the difference between $c_{50}$ values with and without attention ($R_{max}$ ratio: 0.98, $c_{50}$ difference: -6.43) . Right: in the "small attentional field, large stimulus" condition, attention mainly affected response gain, defined as the ratio of $R_{max}$ values ($R_{max}$ ratio: 1.39, $c_{50}$ difference: -0.88). Dotted lines show the percent change in firing caused by attention.

295 sponse gain have also been shown experimentally. Martinez-Trujillo and
296 Treue (2002) recorded from neurons in area MT while presenting two stimuli
297 within the receptive field. One stimulus was moving in a preferred direction,
298 and the other in a non-preferred direction. They then varied the strength of
299 the preferred stimulus while holding the contrast of the non-preferred stim-
300 ulus fixed, and directed the monkey to attend either to the non-preferred
301 stimulus or outside of the receptive field. They found that attending to
302 the non-preferred stimulus caused predominantly a change in contrast-gain.
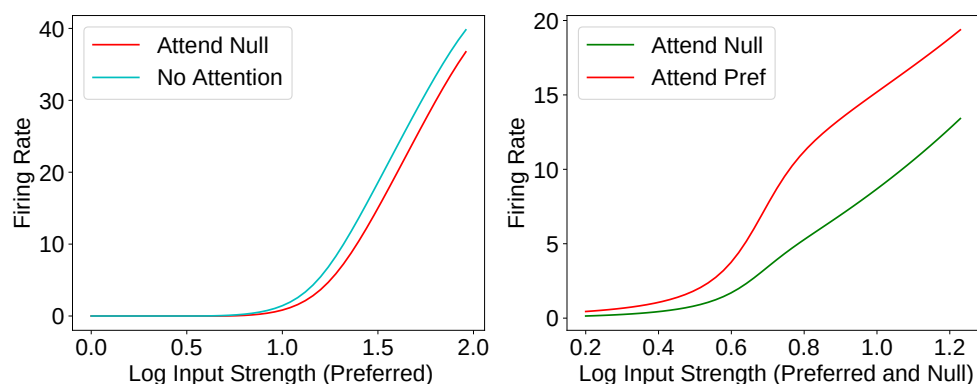
15

Figure 9:

**Experimental paradigm alters gain change type. A.** In the ring model, in the presence of a fixed-strength non-preferred stimulus, the contrast of a preferred stimulus was varied logarithmically (base 10) while attention was directed either away (cyan) or towards the non-preferred stimulus (red) as in Figure 4 of Reynolds and Heeger (2009). Attention to the non-preferred stimulus produced mainly a reduction in contrast gain, measured as the difference between $c_{50}$ values ($R_{max}$ ratio: .97, $c_{50}$ difference: 5.94) (Martinez-Trujillo and Treue, 2002). **B.** Showing preferred and non-preferred stimuli of equal but varying contrast while attending to one or the other produced a much larger change in response gain, measured as the $R_{max}$ ratio ($R_{max}$ ratio: 1.38, $c_{50}$ difference: -2.17). This was studied experimentally in Lee and Maunsell (2009).

However, Lee and Maunsell showed that if the contrast of both the preferred and non-preferred stimulus were varied simultaneously, attending to one or the other stimulus would produce a much larger change in response gain (Lee and Maunsell, 2009). Using the ring model again, we modeled both of these stimulus conditions, and find analogous results (Figure 9A, B).

*2.3.2. Effect on length tuning*

The impact of spatial attention on length tuning was explored in Roberts et al. (2007). In this study, the length of an oriented bar was varied as firing rates from V1 cells were recorded. Attention was directed to the stimulus or to a stimulus in the opposite hemifield. The authors found that, for receptive fields near the fovea, attention had the effect of decreasing preferred length
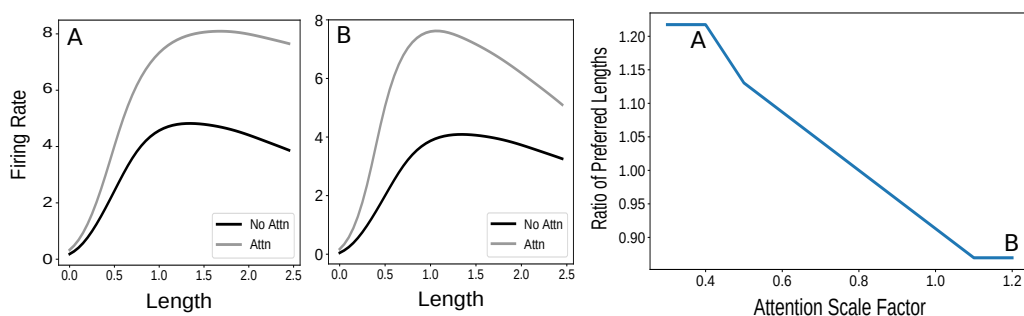
16

Figure 10:

**Size of attention influences length tuning.** Using the line model, we presented a stimulus of increasing length (left two plots). If attention was small compared to the stimulus (far left) attention shifted the preferred length (i.e., the length that elicits the highest firing rate) rightward, making it larger. If the area to which attention was applied was large compared to the stimulus (middle), the opposite occurred. Thus, varying the ratio of the size of attention to the stimulus size ("attention scale factor") caused a shift in the ratio of the preferred lengths (preferred length with attention divided by preferred length without attention; right plot). Scale factor in the far left plot is marked on the right plot by the letter A, middle by B. In Roberts et al. (2007) the ratio of preferred lengths for parafoveal receptive fields was .88 and for peripheral receptive fields 1.19.

314 (that is, the length of the bar that elicits the highest firing rate). For receptive
315 fields in the periphery, the reverse was true: attention increased the preferred
316 length.

317    We explored attention's impact on length tuning using the spatial line
318 model. For different lengths of the stimulus, firing rates were recorded from
319 a neuron at the center. The effect of attention varied as a function of the size
320 of the attentional field. In Figure 10 (right) the ratio of the size of attention to
321 the size of the stimulus is on the x-axis. By keeping a fixed ratio of attention
322 size to stimulus size, we assume that the size of the attentional field scales
323 with the size of the stimulus, but this scaling factor may differ for different
324 cells. For small values of this attention scale factor, the preferred length with
325 attention was greater than the preferred length without it. For higher values,
326 this ratio was reversed. Firing rate as a function of length for two different
327 values of the attention scale factor are shown on the left. This pattern of

17

how attention impacts preferred lengths reflects the impact of attending to the suppressive surround. With attention larger than the stimulus, more of the suppressive surround is activated for any given stimulus length. This effectively increases the length of the stimulus, making the preferred length smaller than without attention.

Our results combined with the findings of Roberts et al. (2007) suggest that attention targets parafoveal receptive fields differently than it targets peripheral ones. In particular, spatial attention inputs to parafoveal cells may be larger than the size of the stimuli these cells respond to. In the periphery, spatial attention inputs may represent an area smaller than the stimuli. This could be a result of the differently sized receptive fields in these two regions.

### 2.3.3. Factors influencing the magnitude of attentional effects

In Lee and Maunsell (2010), the authors controlled attention and task difficulty across stimulus conditions while varying the number of stimuli in the receptive field of MT neurons. Through this, they showed that attentional modulation is weaker when only one stimulus is present in the receptive field, and that this result is well-captured by a divisive normalization model. We use the ring model to replicate these results. By presenting three different stimuli (a most-, moderately-, and least-preferred orientation) either alone or in pairs (Figure 11, left; compare to Lee and Maunsell (2010) Figure 4), we show that the effect of an attentional input was strongest when applied to one stimulus in a pair. In particular, effects of attention on firing rates were highest when moving attention from outside the receptive field to the preferred stimulus inside the receptive field when a non-preferred stimulus is also present (Figure 11, right). The next strongest effect was from moving attention from the non-preferred stimulus in the receptive field to the preferred. Finally, attention effects were weakest when moving attention from outside the receptive field to a preferred stimulus presented alone inside the receptive field.

A similar comparison was done using spatial attention rather than feature attention in Sundberg et al. (2009). Here, attention was moved between the receptive field center and the suppressive surround. A stimulus of preferred orientation was present in the center and was present or absent in the surround. The impact of attending the center was larger when the stimulus in the surround was present (Figure 2 of Sundberg et al. (2009)). We replicated these results using the line model. The firing rate of an excitatory cell was
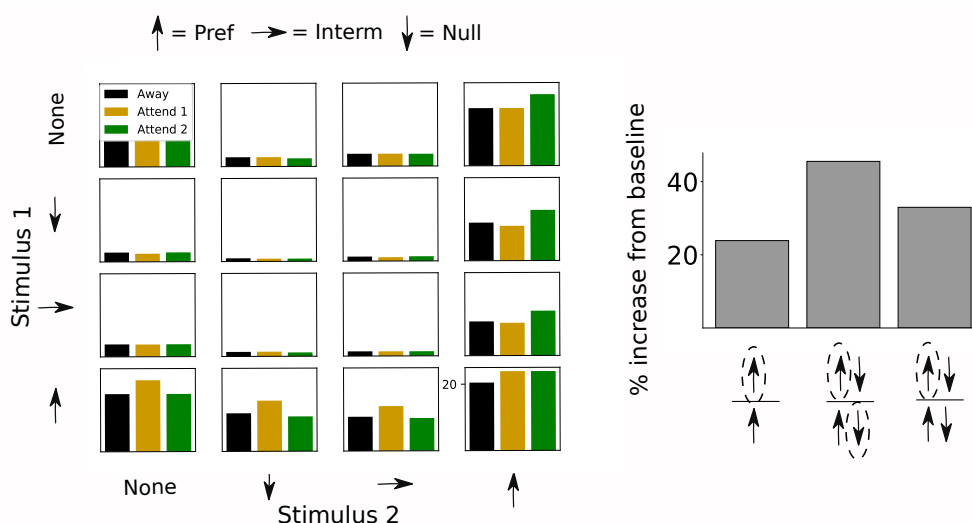
18

Figure 11:

**Effects of attention are greater with more than one stimulus in the receptive field.** Using the ring model, three different stimuli (preferred, intermediate, and null) were shown either individually or in pairs. Attention was directed to either of the two stimuli ('Attend 1' or 'Attend 2') or outside of the receptive field ('Away'; when only one stimulus was present, attending to the opposite stimulus is the same as attending away). Left: Bar plots represent steady state firing of the recorded neuron for all stimulus and attention conditions. Right: bar plots indicate percent increase in firing rate with attention, for three different comparisons. Arrows indicate which stimuli were in the receptive field for the two conditions being compared (bottom arrows indicate baseline condition, top arrow(s) indicate attended condition) and dashed circles indicate attended stimulus. The comparable values for these conditions from Lee and Maunsell (2010) are 9%, 59%, 28% respectively.

365 recorded with a stimulus centered on its preferred location. Attention was
366 applied to this location, or to a location in the surround both in the presence
367 and absence of a stimulus there. There results of this are shown in Figure 12
368 (left).
369     In Sundberg et al. (2009), the impact of attention on surround suppression
370 was also shown over time. The extent to which firing rates are decreased by
371 the presence of the surround was measured when attention was directed to
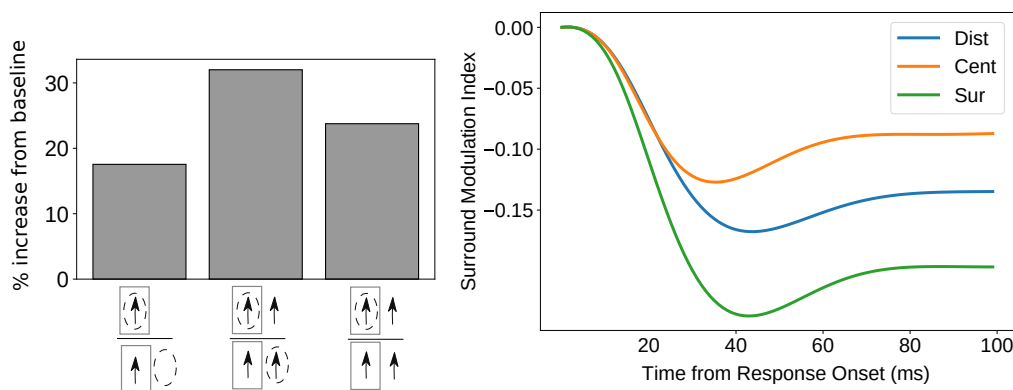
19

Figure 12:

**Effects of attention are greater with a stimulus in the surround.** Using the line model, a preferred stimulus was presented in the receptive field center. Left: bar plot indicates increase in firing in preferred-attended condition (top arrows) vs. baseline condition (bottom arrows). Rectangles indicate receptive field. The presence of a surround stimulus is indicated by an additional arrow outside the receptive field and attention is indicated by a dashed circle. The increase in firing was smaller without the surround present (comparable values from Sundberg et al. (2009) are 18.8% versus 36.8%. The authors do not report the percent increase compared to a baseline condition without attention to either center or surround). Right: the strength of firing rate modulation from the addition of a surround stimulus (the surround modulation index: $[r(C + S) - r(C)]/[r(C + S) + r(C)]$) is plotted vs. time, for different attention conditions: attending the surround, attending the center, and attending a distant location (modeled as no attention). The difference between these conditions emerged over time.

the receptive field center, surround, or to a distant location. The authors note (their Figure 5) that the difference in surround modulation between these different attention conditions emerged over time. The model shows the same result (Figure 12, right). The differences emerge faster in our model than in the data (in the data, the difference is not seen in the time bin 15-55ms after response onset, but emerges sometime in the next 40ms time bin). However, our model does not take into account any delays in the onset of the attentional signal relative to the onset of stimulus-driven feedforward input to the recorded neurons.
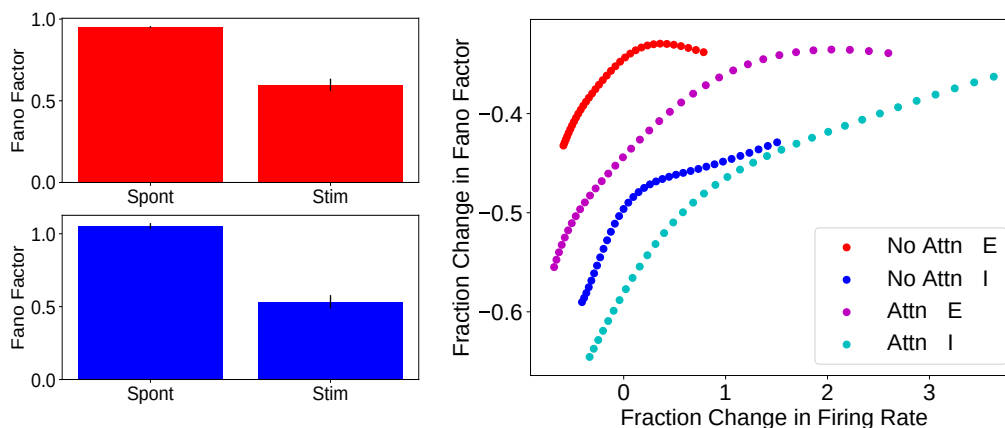
20

Figure 13:

**Attention causes a reduction in trial-to-trial variability.** In the ring model with noisy background input, 35 E (red) and 35 I (blue) cells were recorded as a stimulus that was oblique (but not orthogonal) to their preferred stimuli was presented. Stimulus onset produced a substantial reduction in trial-to-trial variability, measured as the Fano factor, compared to spontaneous activity (left; errorbars are STD). Next, the effect of an attentional modulation was observed. On the right, fractional change in Fano factor is plotted as a function of fractional change in firing rate for each of the 35 E and 35 I cells in the presence and absence of attention. In all cells, stimulus onset produced a decrease in the trial-to-trial variability, regardless of whether the stimulus produced an increase, decrease, or no change in the mean firing rate (Churchland et al., 2010). In the presence of attention, this decrease in variability was enhanced, as has been observed experimentally (Mitchell et al., 2007). The percent change in both firing rate and Fano factor was calculated for each cell by taking a time average of both the mean rate and Fano factor before and after the onset of the stimulus (in trials with attention, it came on at the same time as the stimulus).

<sub>381</sub> *2.4. Attention reduces trial-to-trial variability and noise correlations*

<sub>382</sub> In addition to its effects on mean firing rates, attention has also been
<sub>383</sub> shown to modulate the variability in rates across trials. Mitchell et al. (2007)
<sub>384</sub> showed that attending to a stimulus decreased the across-trial variability of

neural responses when compared to trials in which attention was directed elsewhere. Furthermore, this experiment showed that this decrease in variability occurs in both broad spiking (putative excitatory) cells and narrow spiking (putative inhibitory) cells.

To study this effect in our model, we introduced a source of trial-to-trial variability into our ring network by given each neuron a noisy input in addition to its stimulus inputs, similarly to Hennequin et al. (2018) (see Methods 4.1.1 for details). We then ran 1,000 trials of a simple stimulus presentation. On half of these trials, attention was directed towards the stimulus being presented. On the other half there was no attentional modulation added to the network. The stimulus onset produced a reduction in the trial-to-trial variability, measured as the Fano factor, with this reduction occurring both for neurons that are activated by the stimulus and neurons that are not activated or suppressed (Figure 13), as in experiments (Churchland et al., 2010) and as previously shown for the SSN (Hennequin et al., 2018). Addition of attention caused an additional drop in Fano factor, again regardless of whether the stimulus plus attention caused a net increase, zero change, or net decrease in firing rate (Figure 13, right).

In addition to causing a drop in trial-to-trial variability, Cohen and colleagues demonstrated that an even stronger effect of attention on network variability is a pronounced decrease in the magnitude of noise correlations between neurons in V4 (Cohen and Maunsell, 2009). This aligns with the finding that a stimulus suppresses the shared or correlated component of neural variability, not the component private to each neuron (Churchland et al., 2010). Cohen *et al.*, 2009, recorded from thousands of pairs of neurons and multiunit clusters in V4 during a visual change detection task, and found that the presence of attention greatly enhanced performance. They further showed that the significant improvement in performance was not due to changes in single neurons, but rather to a pronounced drop in the corre-
e correlations). To simulate this experiment, we recorded from pairs of excitatory cells in the ring model in the presence of noisy input while presenting the network with two high-contrast oblique stimuli. On half of the trials, attention was directed to one of the stimuli. We calculated the correlation between all pairs of recorded neurons in the presence and absence of attention. Pairs of neurons were grouped based on their distance from each other on the ring (i.e. difference in preferred orientation). The changes in firing for two example neurons with attention as well as the noise correlations between them over the course of an example trial are shown in Figure 14
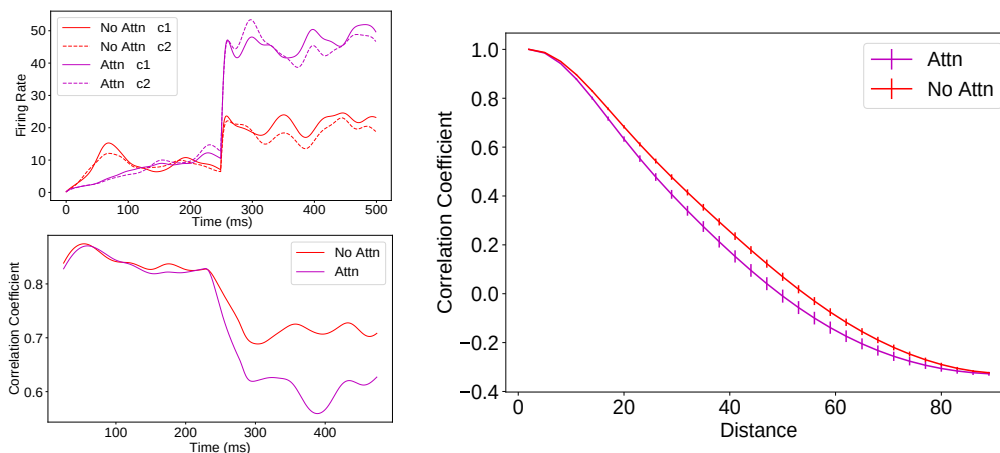
22

Figure 14:

**Attention decreases noise correlations between neurons.** In the ring model with noisy background input, stimulus onset produced a reduction in noise correlations between pairs of neurons in the network. The correlation in firing rates between each pair of cells was calculated as a function of time for each of the two conditions. On the left, an example pair is shown. The mean firing rates of two excitatory cells in each of the two conditions is plotted on top; stimulus (at 90 degrees) and attention turn on at 250ms. The correlations between the two cells are plotted on the bottom. Correlation time-series are shown as a running average with a 50-ms sliding window. On the right, the mean correlation between pairs of recorded cells (representing 30-65 degrees) during the stimulus response epoch is plotted against difference in preferred orientation. Error bars indicate SEM.

423   (left). The average value of noise correlations between neurons at various
424   distances is shown on the right. As was observed experimentally, attention
425   caused a reduction in the noise correlations between neurons beyond the
426   reduction caused by the stimulus alone.

427      The suppression of correlated variability can be understood as resulting
428   from the normalization performed by the model (although it also explains
429   further aspects of this suppression not explained simply by normalization,
430   Hennequin et al. (2018)). In particular, as has been observed experimen-
431   tally (Busse et al., 2009), this normalization averages the responses to ap-
432   proximately equal strength inputs but performs a more unequal averaging

23

of unequal strength stimuli, becoming "winner-take-all" when inputs differ sufficiently in strength (Rubin et al., 2015). The reduction in correlated variability with increasing stimulus strength can be understood to occur because the ongoing noisy inputs become steadily weaker relative to the stimulus. The normalization thus increasingly favors the response to the stimulus and suppresses the noise. Because this suppression is mediated by the network, it acts on the correlated component of the noise and not on the private noise, which is largely averaged out in its impact at the network level.

An alternative picture of the mechanism of suppression is that it occurs through the enhancement of the strength of feedback inhibition with increasing network activation (Hennequin et al., 2018). In particular, in linearizations about the deterministic fixed point, the real parts of the leading eigenvalues become more negative with increasing mean stimulus drive, representing increased feedback inhibition of the corresponding eigenvector activity patterns onto themselves, dampening their fluctuations. Given structured connectivity, these activity patterns have similar structure and so their fluctuations represent correlated variability.

Investigations regarding noise correlations have indicated that a decrease in correlation with attention should only occur for pairs of neurons that represent the same stimulus whereas pairs of neurons representing different spatial locations or features may actually see an increase in correlations (Averbeck et al., 2006). This bi-directional effect of attention was found in area V4 (Ruff and Cohen, 2014). In our ring model, this result occasionally occurred when using weaker stimuli and/or a smaller number of trials to calculate the correlations in the ring model. Examples of this can be found in Supplementary Figure A.17.

The task in Ruff and Cohen (2014), however, used spatial rather than feature attention. Specifically, subjects were required to perform a contrast discrimination task in the cued hemifield. To replicate this study directly we used the line model with two nearby stimuli of unequal contrast (Figure 15, left). The TTS metric from Ruff and Cohen (2014) measures the extent to which a pair of cells have the same (positive TTS) or opposite (negative TTS) preferred stimulus of the two presented. Replicating figure 5 from that paper, we see that attention decreased correlations for cells with the same preferred stimulus but increased it for those with opposite preferred stimuli (Figure 15, right).
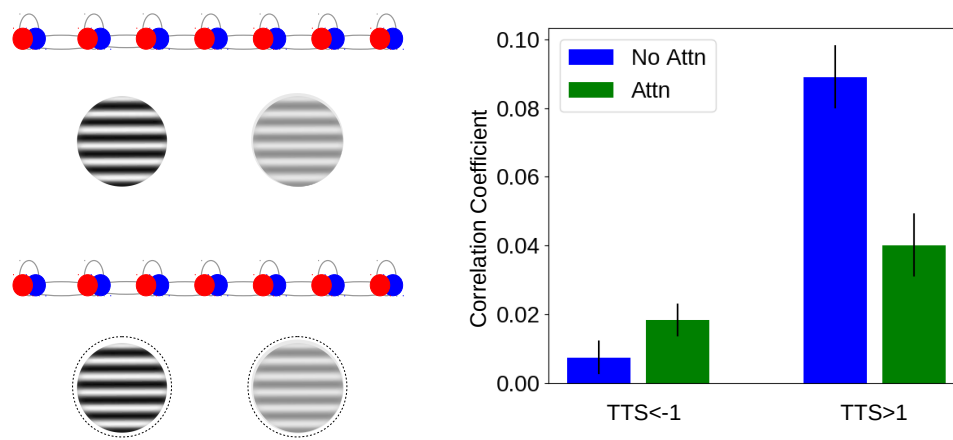
24

Figure 15:

**Attention increases or decreases noise correlations between neurons based on preferred stimulus.** In Ruff and Cohen (2014), animals performed a contrast discrimination task on two nearby stimuli, represented here as two inputs to the line model of different strengths. During different blocks, attention was directed to one of two such sets of stimuli, one in each hemifield. Here we model attention to the opposite hemifield as a 'no attention' condition (top left) and attention to the hemifield of the recorded cells as attention to each of the two stimuli simultaneously (bottom left). The 25 model cells we analyzed responded to one or the other stimulus alone. TTS values are the product of d-primes and represent whether a pair of cells has the same (positive) or different stimulus preference (negative). By creating 20 populations of 25 cells each, we analyzed the relationship between TTS and the effect of correlation on attention for 6000 cell pairs. Through this we found both a significant ($p << .05$) decrease in correlation with attention for cells that preferred the same stimulus and increase for cells that had opposite preferences (right). Error bars indicate SEM. For more details, see Methods 4.2.

### 2.5. An alternative mechanism

In all of the simulation results presented thus far, attentional modulation has been modeled as a small excitatory input biased towards the excitatory cells within the locus of attention. Here we consider instead a small in-

25

hibitory input to inhibitory cells within the locus of attention, disinhibiting rather than exciting the excitatory cells. This is motivated by two observations. First, it was observed that inputs from Anterior Cingulate Cortex to V1 target the VIP class of inhibitory cells (Zhang et al., 2014). The VIP cells in turn are known to inhibit other inhibitory neurons and, at least in V1, disinhibit excitatory cells (*e.g.* Fu et al., 2014). The ACC input conceivably could be involved in attentional modulation. Second, recent electrophysiologic work has revealed the function of two classes of inhibitory cells in layer 1 of cortex (Jiang et al., 2013). One of these classes, the single bouquet cells (SBCs) was shown to preferentially inhibit the interneurons of deeper layers, and so have a net disinhibitory effect on the local pyramidal cells. As layer 1 receives a significant portion of its input from higher cortical areas, it has been suggested that this circuit may play a role in attention and other top-down modulation of local circuit activity (Larkum, 2013).

To test the feasibility of this mechanism in our model, we repeated our suite of simulations using this alternative, disinhibitory mechanism of attention. Rather than modeling attention as an additional excitatory input to E cells, we instead model it as an additional inhibitory input to I cells. The results of these simulations are presented in the Supplementary Figures. Overall, this alternative mechanism can qualitatively reproduce most of the findings we report above (Supplementary Figure A.18). Frequently, however, the same value of the attention strength parameter produces weaker effects on neural firing than when attention is directed towards the excitatory cells (for example, compare Figure 12 to Figure A.18G).

In addition, there are instances where this form of attention does not qualitatively replicate our original findings (Supplementary Figure A.19). One major discrepancy between results comes from the use of the 2-D model. Comparing Supplementary Figure A.19B to Figure 5, modeling attention as inhibition to inhibitory cells creates the opposite relationship (i.e., a negative correlation) between attentional modulation and normalization. In the 2-D model, any additional inhibitory input to the inhibitory population has the effect of increasing firing rates for many of the cells, even those representing unattended stimuli. The model therefore cannot replicate findings that rely on attention to a non-preferred stimulus causing a decrease in firing rate. This appears to be a consequence of the strong inhibition needed to keep this more complex model in a stable regime. Attention directed toward inhibitory cells also has a surprising effect on the correlations explored in Figure 15. As can be seen in Supplementary Figure A.19E, this fom of attention increases

correlations for pairs of cells both with the same and opposite preferred stimuli.

## 2.6. Attention enhances detection performance in a multi-layer model

An important consequence of deploying attention is enhanced performance on challenging tasks. We have thus far shown how the SSN can replicate many neural effects of attention, but to truly understand attention, it is necessary to link these neural changes to performance changes. And for that it is necessary to build a functioning model of the visual system that can perform visual tasks.

Because the SSN replicates neural findings that have been found in various areas in the visual system, it can be thought of as a canonical circuit, which is repeated throughout the visual hierarchy. To build a biologically-realistic multi-area model of the visual system that can perform a task, we model each area as a set of SSNs, the outputs of which are fed into another set of SSNs (i.e., a downstream visual area). The precise connections between these areas are learned as part of a training procedure. In particular, the SSN circuitry is placed inside a convolutional neural network architecture, creating a model we have dubbed the SSN-CNN (Methods 4.3).

The structure of the model can be seen in Figure 16A. The network is a 2-layer convolutional neural network wherein the convolutional filters are constrained to be non-negative (to mimic the excitatory feedforward connections that exist between different visual areas). In addition, after each pooling layer is an SSN layer. The SSN layer implements normalization (historically normalization layers have been included in CNNs, typically implemented via a divisive normalization equation Krizhevsky et al. (2012)). Specifically, at each 2-D spatial location, a ring SSN implements feature normalization across the different feature maps. The recurrent connections of the SSN layers are held constant while all other weights of the network are trained end-to-end via backpropagation through time on the MNIST 10-way digit classification task.

After the network is trained on the standard task, the final layer is replaced by a series of binary classifiers, one for each digit. These binary classifiers are trained on digit images to determine if a given digit is present in the image or not (for example, one of the binary classifiers would be trained to classify images as being of the digit '4' or not). To test the impact of attention on the abilities of these binary classifiers, we presented the network with a more challenging task: determining if a given digit is present in an image
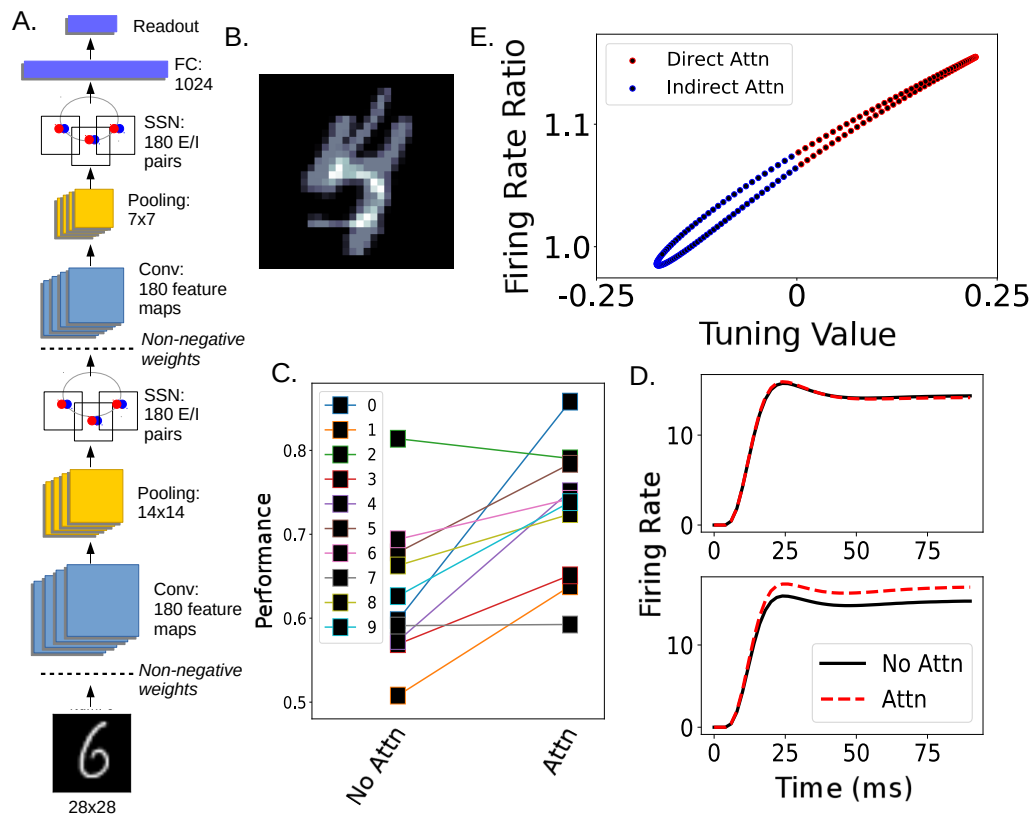
27

Figure 16:

**Attention in the SSN-CNN enhances visual detection performance.** A.) The architecture of the SSN-CNN model. In the SSN layers, a full ring model exists at each spatial location (though only one is shown). B.) An example of the images used in the attention task. This image contains a '5' and '4' overlaid, therefore both the binary classifier trained to detect 4s and the one trained to detect 5s should respond positively. C.) Binary detection performance for each digit with (right) and without (left) attention. D.) Example firing rate of two neurons recorded from the second SSN layer with receptive fields at the center of the image when shown the image in (B). The top neuron had a small decrease in firing when attention was deployed to the digit 4 and the bottom had an increase. E.) Impact of attention to the digit 4 on firing rates of excitatory cells (rate with attention divided by rate without) as a function of tuning to the digit. A feature map's tuning value for a given digit is defined as its z-scored mean response to that digit (see Methods, section 4.3). Attention is modeled as excitatory input applied to feature maps whose tuning value is above the median value across maps for that digit. The strength of a map's attentional input is proportional to the difference between that map's tuning value and the median value. Only neurons marked in red were above the median and given direct attentional input.

that contains two overlaid digits (Figure 16B). The network performs above chance on this challenging task, and performance increased when attention was applied (Figure 16C, attention applied at layer 2).

Attention is applied in this model as previously described: an additional positive input is given to excitatory cells that prefer the attended digit. To determine which cells in the SSN layers "prefer" the attended digit we created tuning curves based on the response of excitatory cells in the SSN when presented with images of different digits (See Methods 4.3). Applying attention in this way still elicits attentional changes in the cells that are not directly targeted—through the recurrent connections—as can be seen in Figure 16E. This includes decreasing the firing rates of neurons that do not prefer the attended digit. While this feature attention is applied the same way across all ring networks at a layer, the pattern of feedforward input will influence the ultimate impact of attention. This can be seen by comparing the ratio of firing with and without attention in ring networks at different nearby spatial locations, which receive slightly different feedforward input (Supplementary Figure A.20).

Previous work (Lindsay and Miller (2018); Lindsay (2015)) has shown how attentional changes in different layers of a deep convolutional neural network can lead to enhanced performance on challenging visual tasks. That work demonstrated that the attentional modulation style that works best is multiplicative and bi-directional changes (i.e., the effect of attention should be to scale the activity of neurons that prefer the attended stimulus up and those that don't prefer it down). What we have shown here is how an additive input solely to the excitatory neurons that prefer the attended stimulus can turn into multiplicative and bi-directional changes via the circuit mechanisms of the SSN and lead to an increase in performance. This allows for a straightforward mechanism by which top-down attentional signals can lead to enhanced performance simply by providing additional synaptic inputs to the right set of excitatory cells.

## 3. Discussion

The stabilized supralinear network (SSN) is a model of recurrent processing in visual cortex that is informed by anatomy and replicates several features of neural activity (Rubin et al., 2015). With a simple addition to this nonlinear circuit model, we are able to reproduce a number of experimental results on attention in visual cortex (Treue and Martinez Trujillo,

29

1999; Cohen and Maunsell, 2009; Mitchell et al., 2007; Reynolds and Desimone, 2003; Sundberg et al., 2009; Lee and Maunsell, 2009; Ni et al., 2012; Martinez-Trujillo and Treue, 2002). Through balanced amplification (Murphy and Miller, 2009), a small additional excitatory input to excitatory cells causes a nonlinear scaling of firing rates in a manner consistent with a number of experimental observations. Recurrent connections implement interactions between features and spatial locations. These simple models are able to account for changes in stimulus interactions, differences in gain changes and the magnitude of attention's effects, as well as changes in trial-to-trial variability. We are not aware of any previous model that has attempted to replicate so many effects of attention simultaneously. The ability to replicate all these effects via a small additional input to a subset of neurons provides a simple, plausible mechanism through which higher cortical feedback can implement attention.

Previous work has identified areas in the frontal cortex that may be considered the source of top-down selective visual attention (Bichot et al., 2015; Paneri and Gregoriou, 2017). Exactly how connections from these areas target visual areas to create the changes seen with attention is unknown. Studying these feedback connection can be challenging, as it requires detailed anatomical investigations across multiple brain areas. For this reason, narrowing the hypothesis space by identifying which mechanisms of feedback control are theoretically capable of implementing the known effects of attention is important. Here, we show that positive additive input to the excitatory neurons that prefer the attended stimulus can recreate the multiplicative changes observed in both E and I cells and both in cells that prefer and do not prefer the attended stimulus. Adding negative input to the inhibitory cells that prefer the attended stimulus can also replicate most of these effects, except that in our 2D model it tended to raise firing rates of neurons that did not prefer the attended stimulus. We do not know if that is a fundamental problem with a disinhibitory model of attention or if it could be fixed by altering model connectivity. Overall, these results show that feedback connections do not need to be directly responsible for all of the neural effects of attention. Instead, they only need to target a subset of neurons in a simple specific way and the local recurrent circuitry can take care of the rest.

There are effects of attention that this model does not readily replicate. For example, spatial attention has been observed to shift and shrink receptive fields. A previous two-layer model with multiplicative attentional inputs

and inhibitory recurrent connections was able to replicate these phenomena (Miconi and VanRullen, 2016). Creating a unified model that can capture all of attention's relevant effects is a goal for future work.

In addition to replicating known findings, the set of models presented here can serve as testbeds for future work on attention. In particular, experimental designs can be explored and precise predictions made before carrying out further experiments.

Circuit models in neuroscience are frequently built to replicate and understand the relationship between anatomy and neural activity. Traditionally, these models do not perform a perceptual or cognitive task. Yet, an ultimate understanding of the circuitry of visual perception will need to replicate behavioral as well as neural findings. We work towards this goal here by incorporating the SSN model into a convolutional neural network that can perform digit recognition (the SSN-CNN). Through this, we connected the neural changes our model replicates to enhanced detection performance. This model also sets a precedent for how traditional approaches from computational neuroscience can be incorporated with the increasingly popular approach of using deep neural networks to study the brain (Yamins and DiCarlo, 2016; Kell and McDermott, 2019).

Further connections between neural changes and performance remain to be explored, and the SSN-CNN could be useful in this pursuit. For example, we do not incorporate noise into the SSN-CNN in this work, however using the noisy version of the ring model (Figures 13 and 14) would allow for an exploration of how noise and correlation changes impact performance. We also do not attempt to model or replicate effects of attention on reaction time, however that is possible in this dynamic model. Using the full 2-D model (instead of ring models at each spatial location) would also allow for an exploration of the effects of spatial attention and the interaction between spatial and feature attention.

## 4. Methods

Code will be publicly available upon publication.

### 4.1. Basic Circuit Models

In this study we employ several different configurations of a basic SSN circuit model, the central unit in all being an interconnected pair of excitatory (E) and inhibitory (I) cells. The two core models are the one-dimensional

ring model and the one-dimensional line model. In addition, for Figure 1 we use a simplified 2-cell circuit model, and for Figures 5 and 7 we use a large two-dimensional model.

In all models, each neuron, $i$, is represented as a firing rate unit whose activity, $r_i$, evolves according to:

$$\tau_i \frac{d}{dt} r_i = -r_i + k \left( [I_i]_+ \right)^n \tag{2}$$

with $n > 1$ (indicating a supralinear activation function). The expression $[v]_+ = \max(v, 0)$, that is, neuronal activity cannot go below zero. The inputs, $I_i$, to a given neuron $i$ are comprised of recurrent inputs, feedforward stimulus inputs, and attentional inputs. These inputs and parameter values are specified for each model below. In all models the time constant $\tau_i$ has the value $\tau_E = 20$ m$s$ for all E cells and $\tau_I = 10$ m$s$ for all I cells. Simulations are run using the forward Euler method with time step 1ms.

All of these models except the E-I pair model were described previously in (Rubin et al., 2015), however we will recap them briefly here. We used all the same model parameters from that study, and did not tune them in any way to get the current results. The models are only modified by the addition of attentional inputs, and by the addition of noise inputs for Figures 13 and 14.

### 4.1.1. Ring Model

The ring model is intended to represent neurons with a shared retinotopic receptive field but different preferred features. In this model, an E-I pair exists at each location on the ring, with the preferred feature (e.g. orientation or direction) varying smoothly around the ring. The relative input to a cell with preferred orientation $\theta$ from a stimulus of orientation $\phi$ is given by $h(\theta, \phi) = e^{-\frac{d_{circ}(\theta - \phi)^2}{2\sigma_{FF}^2}}$ where $d_{circ}(\theta - \phi)$ is the shortest distance around the circle between $\theta$ and $\phi$. The absolute stimulus input to a cell comes from multiplying $h(\theta, \phi)$ by the scalar $c$, which represents the overall strength or contrast of the stimulus. In addition, attention directed towards orientation $\phi'$ provides extra input to E cells with the same overall shape as a stimulus input, scaled by the attention strength factor, $a$. (In studies that modeled attention as negative input to I cells rather than positive input to E cells, this input is instead given to inhibitory cells, with $a < 0$.) In total, input to

the E or I cell at location $\theta$ on the ring is given by:

$$
\begin{aligned}
I_E(\theta) &= ch(\theta,\phi) + ah(\theta,\phi') + \sum_{\theta'} W_{EE}(\theta,\theta')r_E(\theta') - W_{EI}(\theta,\theta')r_I(\theta') \\
I_I(\theta) &= ch(\theta,\phi) + \sum_{\theta'} W_{IE}(\theta,\theta')r_E(\theta') - W_{II}(\theta,\theta')r_I(\theta')
\end{aligned}
\tag{3}
$$

respectively.

Recurrent connections fall off according to $W_{ab}(\theta - \theta') = J_{ab}e^{-\frac{d_{circ}(\theta-\theta')}{2\sigma_{ori}^2}}$, where $d_{circ}(\theta - \theta')$ is the shortest distance around the circle between $\theta$ and $\theta'$. If multiple stimuli are present the inputs are added linearly.

For simulations of this model, the following parameters are used: the number of E/I pairs is $N = 180$; the spacing in degrees between adjacent pairs on the ring is $\Delta\theta = 1°$; $J_{EE} = 0.044$, $J_{IE} = 0.042$, $J_{EI} = 0.023$, $J_{II} = 0.018$, $\sigma_{ori} = 32°$, $\sigma_{FF} = 30°$, $k = 0.04$, $n = 2.0$.

The ring and its inputs are schematized in Figure 2.

In certain simulations, noise is added to the inputs to these cells. Specifically, $10 + \nu(\theta, t)$ was added to input to each unit at each timestep. External noise $\nu$ was given by convolution of unit-integral Gaussian temporal filter (stdev 10 ms) and spatial filter (stdev 8°) with Gaussian spatiotemporally white noise (mean 0, stdev 40), yielding $\sqrt{\langle \nu^2 \rangle} \approx 1$.

### 4.1.2. Line Model

In the line model, each E-I pair represents a different retinotopic location but all have the same preferred features. Rather than being arranged in a ring, these pairs are simply placed on a line. The line model follows the same basic equations as the ring model, however the stimulus input is defined differently and the recurrent connections are differently arranged.

A stimulus input is defined in terms of stimulus center $x_0$ (taken as zero for center stimuli), length $l$ and sharpness parameter $\sigma_{RF}$. The input to an E-I pair at location $x$ is given by $s_l(x-x_0) = \left( \frac{1}{1+e^{-\frac{(x-x_0)+l/2}{\sigma_{RF}}}} \right) \left( 1 - \frac{1}{1+e^{-\frac{(x-x_0)-l/2}{\sigma_{RF}}}} \right)$. As in the ring model, this input is scaled by the overall strength of the stimulus, $c$.

In this model, there are $N$ E/I units with grid spacing $\Delta x$. Recurrent connections are defined with respect to distance between neurons. Excitatory projections are given by $W_{aE}(x, x') = J_{aE}e^{-\frac{|x-x'|^2}{2\sigma_{aE}^2}}$ for $a \in \{E, I\}$. Inhibitory projections $W_{aI}$ are only to the same line position as the projecting neuron.

33

719      The parameters used in this model are: $N = 101$, $\Delta x = \frac{1}{3}^{\mathrm{o}}$, $\sigma_{RF} =$
720   $0.125\Delta x$, $J_{EE} = 1.0$, $J_{IE} = 1.25$, $W_{EI} = 1.0$, $W_{II} = 0.75$, $\sigma_{EE} = \frac{2}{3}^{\mathrm{o}}$,
721   $\sigma_{IE} = \frac{4}{3}^{\mathrm{o}}$, $k = 0.01$, $n = 2.2$.

722      Again, if multiple stimuli are present their inputs are simply added to-
723 gether and attention takes the same shape as a stimulus but is only directed
724 toward E cells.

725      In one simulation, noise was added to the line model. This noise was
726 similar to that added to the ring model, but with a lower baseline (5 instead
727 of 10) and different spatiotemporal parameters: external noise was given
728 by convolution of unit-integral Gaussian temporal filter (stdev 15 ms) and
729 spatial filter (stdev $3\Delta x$) with Gaussian spatiotemporally white noise (mean
730 0, stdev 10).

731 *4.1.3. 2-D Model*

732      The one-dimensional ring and line models vary either in preferred retino-
733 topic location or visual feature. To create a model wherein cells have both
734 varying retinotopic as well as feature preferences, we place E-I pairs on a
735 two-dimensional spatial grid representing retinotopy, with an overlaid map
736 of preferred orientation (which may be imagined to represent any circular
737 preferred feature). This model also incorporates randomness in parameters,
738 allowing study of diversity in responses as in Fig. 5.

739      Let $W_{ab}(x, x')$ be the synaptic weight from the cell of type $b$ (E or I), at
740 position $x'$, with preferred orientation $\theta(x')$, to the cell of type $a$, at position $x$,
741 with preferred orientation $\theta(x)$. Nonzero connections are sparse and chosen
742 randomly, with probability $p\left(W_{ab}(x, x') \neq 0\right) = \kappa_b e^{-\frac{(x-x')^2}{2\sigma_{ab}^2}} e^{-\frac{d_{circ}(\theta(x)-\theta(x'))^2}{2\sigma_{ori}^2}}$.
743 Where a nonzero connection exists, $W_{ab}(x, x')$ is chosen randomly from a
744 Gaussian distribution with mean $J_{ab}$ and standard deviation $0.25J_{ab}$; weights
745 of opposite sign to $J_{ab}$ are set to zero. For each cell, the set of recurrent
746 synaptic weights of type $b$ (E or I) it receives are then scaled so that all
747 cells of a given type $a$ (E or I) receive the same total inhibitory and the
748 same total excitatory synaptic weight from the network, equal to $J_{ab}$ times
749 the mean number of connections received under $p\left(W_{ab}(x, x') \neq 0\right)$. $\tau_E$, $\tau_I$,
750 $n_E$, $n_I$, and $k$ are also drawn from Gaussian distributions, with standard
751 deviation 0.05 times the mean (parameter values below indicate means).

752      We use a grid of $75 \times 75$ E-I pairs. The preferred orientation of an E-I pair
753 is given by a map randomly generated using the method of Ref. (Kaschube
754 et al., 2010), (their supplemental materials, Eq. 20) with $n = 30$ and $k_c =$

755    $\frac{8 \text{ cycles}}{75 \text{ grid intervals}}$. The full map is taken to be $16° \times 16°$; the grid interval

756    $\Delta x = \frac{16°}{75}$. Boundaries in retinotopic space are periodic. Parameters: $\kappa_E =$

757    $0.1$, $\kappa_I = 0.5$, $J_{EE} = 0.10$, $J_{IE} = 0.38$, $J_{EI} = 0.089$, $J_{II} = 0.096$, $k = 0.012$,

758    $n_E = 2.0$, $n_I = 2.2$, $\sigma_{EE} = 8\Delta x$, $\sigma_{IE} = 12\Delta x$, $\sigma_{EI} = \sigma_{II} = 4\Delta x$, $\sigma_{ori} = 45°$,

759    $\sigma_{FF} = 32°$, $\sigma_{RF} = \Delta x$. Degrees can be converted to distance across cortex

760    by assuming a cortical magnification factor of 0.6 mm/deg, a typical figure

761    for $5 - 10°$ eccentricity in the cat (Albus, 1975) giving $\sigma_{EE} = \sigma_{IE} = 1.54$mm,

762    $\sigma_{EI} = \sigma_{II} = 0.513$mm, orientation map period 1.2mm.

763    In this model, the relative input to the cell at 2D-position $\mathbf{x}$ with pre-

764    ferred orientation $\theta(\mathbf{x})$ from a grating of size $l$ centered at position $\mathbf{x}'$ with

765    orientation $\phi$ is $h(\mathbf{x}) = s_l(|\mathbf{x} - \mathbf{x}'|)e^{-\frac{d_{circ}(\theta(\mathbf{x})-\phi)^2}{2\sigma_{FF}^2}}$; for a full-field grating, the

766    relative input is simply $h(\mathbf{x}) = e^{-\frac{d_{circ}(\theta(\mathbf{x})-\phi)^2}{2\sigma_{FF}^2}}$.

767    We used different exponents, $n_I > n_E$, to increase stability despite vari-

768    ability (as supported by experiments: Supplemental Figure S3 of Ref. Haider

769    et al., 2010). Variability of $\tau$'s, $n$'s, $k$ was limited because larger variabil-

770    ity tended to yield instability; biologically, large variability can probably be

771    tolerated without instability because of various forms of homeostatic com-

772    pensation (Turrigiano, 2011), not modeled here.

773    *4.1.4. E-I Pair Model*

774    In Figure 1 we study an isolated E-I pair. The inputs in this simple

775    two-neuron model are given by:

$$\begin{aligned} I_E &= W_{EE}r_E - W_{EI} * r_I + c_E \\ I_I &= W_{IE}r_E - W_{II} * r_I + c_I \end{aligned} \tag{4}$$

776    We use the following parameters: $W_{EE} = 1.00$, $W_{IE} = 1.25$, $W_{EI} = 0.75$,

777    $W_{II} = 0.75$, $k = 0.01$, and $n = 2.2$. The inputs $c_E$ and $c_I$ are the sums of

778    two components, an "orientation tuned" input that is equal between the two

779    neurons and an untuned modulatory component added to either the E or I

780    cell on a given trial. The tuned component is given by a Gaussian curve at

781    orientation $\theta$: $50e^{-\frac{\theta^2}{2\sigma^2}}$, $\sigma = 20°$. Modulatory input: to I cells, from 0 to 10

782    in steps of 2.5; to E cells, from 0 to 5 in steps of 1.25.

783    *4.2. Attention Experiments*

784    Unless otherwise noted, simulations ran for 300ms and final firing rates for

785    excitatory cells were reported. Attention was modeled as additional input

35

of a specified strength given only to the excitatory cell in a pair. Unless otherwise stated, the shape of the attentional inputs was the same as that of the attended stimulus (as schematized in Figure 4.1.1).

### 4.2.1. Using the Ring Model

In Figure 3, we used the ring model to show how attention to a non-preferred stimulus enhances suppression. The preferred stimulus was oriented at 45 degrees, with strength 40. The non-preferred was oriented at 135 degrees and the strength varied from 0 to 80. Attention was applied to either stimulus at strength 3.

In Figure 4, a non-preferred stimulus (oriented at 135 degrees with strength 40) for the recorded cell (located at 45 degrees) was present as another stimulus (also strength 40) varied from orientation 0 to 180 degrees. Attention (strength 2) was applied to the non-preferred probe stimulus, to the varying stimulus, or not applied at all.

In Figure 9 (left), activity was recorded from a cell at 45 degrees while a preferred stimulus (45 degrees) was presented in conjunction with a non-preferred (135 degrees) stimulus. While the non-preferred stimulus remained at strength 50, the strength of the preferred one varied logarithmically from 1 to 100. Attention was directed to the non-preferred stimulus with strength 5 (or was absent). In Figure 9(right), the contrast of both the preferred and non-preferred stimulus varied logarithmically from ≈1-20. Attention was applied either to the preferred or non-preferred stimulus with strength 1.

In Figure 11, the cell located at 10 degrees was recorded. Each combination of a preferred stimulus (20 degrees), intermediate stimulus (60 degrees), non-preferred stimulus (80 degrees), or no stimulus was tested. All stimuli were presented with strength 20 and an additional input of 10 was given to all cells to better match the baseline firing in (Sundberg et al., 2009). Attention (of strength 1.5) was applied to either of the stimuli present or not at all.

In Figures 13 and 14, the ring model with added noise was used and simulations ran for 500ms. In Figure 13, for the first 250ms, no stimulus or attentional inputs are given (noise inputs are on throughout). At 250ms, a stimulus of strength 25 located at 90 degrees turns on, and on half of the trials so does an attentional input at the same location (strength 8). 1000 trials are run in total. To calculate spontaneous firing rates and Fano factor (FF), firing rates are averaged over 100-250ms. For stimulus-evoked activity, they are averaged over 350-500ms (these are the two epochs compared when calculating the fraction change in firing and FF in the right plot of the figure).

36

Both E and I cells from 30-65 degrees were recorded.

In Figure 14, for the first 250ms, no stimulus or attentional inputs are given (noise inputs are on throughout). At 250ms, two stimuli (both of strength 25, one located at 90 degrees and one at 45) turn on, and on half of the trials so does an attentional input at 90 degrees (strength 8). 1000 trials are run in total. For the figure on the left, correlations are calculated in overlapping windows of 50ms. On the right, correlations are calculated from firing rates averaged over 350-500ms. E cells at all locations were recorded and correlation is plotted as a function of the distance on the ring between any two pairs.

### 4.2.2. Using the Line Model

In Figure 6, a stimulus of strength 25 and length $\frac{14}{15}$ spatial degrees was either placed at the center of the receptive field of the cell at position 0, placed in its surround (at a distance of $\frac{21}{15}$ degrees), or placed at both locations simultaneously. In the last configuration, attention (strength 2) was applied either to the stimulus at the center or the surround (or not at all).

In Figure 8 (left), a stimulus of length 1 spatial degree is presented at the center of the recorded cell with contrast varying logarithmically from 1-100. Attention of strength 1 and length 25 degrees is applied at the same location. For the figure on the right, the size of the attention and stimulus are reversed. To replicate differences in baseline firing shown in (Reynolds and Heeger, 2009), an additional input of 10 is given to all cells in the simulations producing the figure on the left, and an additional input of 2 is given for those on the right.

In Figure 10, a stimulus of strength 15 was centered on the receptive field of the recorded cell with length varying from 0 to 2.5 degrees. The size of attention (applied with strength 4) was equal to the length of the stimulus times an attention scale factor which ranged from .3 to 1.2. The preferred length is defined as the length at which the maximal firing rate is elicited.

In Figure 12, a stimulus of length 1 degree and strength 25 is centered on the recorded neuron's receptive field. A stimulus of the same size and strength either is or isn't presented in the surround (1.5 degrees away). Attention (strength 1, length 1) is applied to the center or surround location in each condition.

In Figure 15, the line model with noise added is used. Two stimuli each of length 2.75 degrees were placed at a distance of 2 degrees on either side of the center of the line model. One had a $c$ of 30 and the other 65. On

attention trials, attention was applied to both stimuli with a strength of 5. For each 'recording session' simulated, excitatory cells 39-63 (roughly 4 degrees on either side of the center cell at 51) were recorded as these cells responded to one or the other stimulus alone. Responses to each stimulus alone at $c = 65$ (50 trials each) were used to calculate a d-prime value for each cell that represents the extent to which that cell prefers one stimulus over the other. As in Ruff and Cohen (2014), the product of d-primes defined the TTS (task tuning similarity) value for a pair of cells. 100 attention trials and 100 no attention trials were run to calculate the correlation coefficients for each pair of cells in each condition based on the average firing over the final 25ms of the simulation (results are the same using 250 or 500 trials). 20 different 'recording sessions' were created using a different random seed for the noise with each one. In addition to the mean changes plotted in Figure 15, we also explored the relationship between TTS and correlation by fitting separate lines to the correlation versus TTS plot in the no attention case and the attention case. If attention differently affects negative and positive TTS pairs, the slope of the attention line should be less than the no attention line. Using the same bootstrap analysis as in Ruff and Cohen (2014) we found this to be true for all 20 of our populations (not shown).

### 4.2.3. Using the 2-D Model

In Figure 5, the two-dimensional model was used to explore the relationship between normalization and attention. We sampled 250 excitatory cells from the model. For each cell, a stimulus of preferred orientation, size 16 degrees, and strength 40 is presented to the cell. An orthogonal stimulus of the same size, position, and strength (the "null" stimulus) is then presented, and then the preferred and orthogonal stimuli are presented together. Attention (strength 8) is applied either to the preferred or null stimulus. These response values are used to calculate the normalization modulation index and attention modulation index for each cell.

In Figure 7, we sample 100 cells from the model to test the interaction between surround suppression and attention. For each cell, a stimulus of strength of 50 of preferred orientation and size 10 degrees is shown. A stimulus with the same orientation and strength is placed in the surround at a distance of 10 degrees, and the response is recorded. The surround at 10 degrees is, technically, a circumference of possible positions around the center. To decide where to place the surround stimulus, the surrounding neuron at a distance of 10 with a preferred orientation closest to that of the center neuron

897  is chosen. Attention (modulation strength = 5) is then directed either to the
898  center or surround stimulus.

### 4.3. The SSN-CNN Model and Experiments

900  The SSN-CNN is an adaptation of a traditional convolutional neural net-
901  work. The inputs to the network are grayscale images of handwritten digits
902  (28-by-28 pixels). The first convolutional layer applies 180 separate 3×3 fil-
903  ters, all of which are constrained during training to contain only non-negative
904  values. The application of these filters results in 180 feature maps, each with
905  a spatial dimension of 28×28. A 3×3 max-pooling layer with stride 2×2 re-
906  duces the feature map size down to 14×14. The output of the pooling layer
907  determines the input to the ring SSNs that exist at the next layer. Specifi-
908  cally, at each of the locations on the 14x14 spatial map, there is a ring SSN
909  with 180 E/I pairs. The activity of the units in the 180 feature maps provide
910  the $c$ values (that is, the strength) for inputs centered at that location on
911  the ring. We arbitrarily number the feature maps from 1 to 180 and let $\phi$
912  be the number of a particular feature map. Then at spatial position $x, y$,
913  the feedforward input to each cell in the E-I pair located at position $\theta$ in the
914  ring model is given by $\sum_{\phi} c_{x,y}(\phi) \, h(\theta, \phi)$, with $c_{x,y}(\phi)$ the activity of the unit
915  in the $\phi$ feature map in the pooling layer at location $x, y$, and $h(\theta, \phi)$ the
916  function defined in section 4.1.1. While there is no concept of a ring in the
917  topology of the feature maps prior to learning, we still map the 180 feature
918  maps onto the 180 locations in the ring. Because feature maps assigned to
919  more nearby locations in the ring will more strongly influence one another's
920  output on the ring, the feature maps should ultimately develop structure
921  reflecting the ring topology (Lindsay and Miller, 2018).

922  This architecture is then repeated to create a two-layer convolutional
923  network. The output of the second SSN layer serves as input to a fully-
924  connected layer with 1024 units, which then projects to the final 10-unit
925  layer (one for each digit). For training, the network was unrolled for 46
926  timesteps (with $dt = 2$ms for the SSN layers) and trained on the MNIST
927  dataset using backpropagation through time to minimize a cross entropy loss
928  function (batch size 128). Only the final timestep was used for calculating
929  the loss function and classification accuracy. The recurrent weights for each
930  ring SSN at both layers were set as described above for the standard ring
931  network. These weights were not allowed to change during training.

932  Repeating the procedure of (Lindsay and Miller, 2018), once the network
933  was trained on the standard classification task, the final 10-unit layer was

replaced with a series of binary classifiers, one for each digit. The weights from the 1024-unit second-to-last layer to the 2-unit final layer were trained to perform binary classification on a balanced training set wherein half of the images were of the given digit and half without.

We then generate more challenging images on which to test the benefits of attention. These images consist of two regular MNIST images added together. The test set for each binary classifier contains 768 images, half of which contain (as one of the two digits) the digit the classifier was trained to detect and the other half do not. Performance accuracy is given as the overall percent correct of the binary classifier on this test set.

To know how to apply attention, we first present 45 standard MNIST images of each digit to the network and record the activity of neurons in the SSN. From this we calculate "tuning values" that indicate the extent to which each feature map prefers each digit. As in (Lindsay and Miller, 2018), tuning values are defined as a z-scored measure of the feature map's mean response to each digit. Specifically, for feature map $\theta$ in the $l^{th}$ layer, we define $r^l(\theta, n)$ as the activity in response to image $n$, averaged over all units in the feature map (i.e., over the spatial dimensions). Averaging these values over all images in the training sets ($N_d = 45$ images per digits, 10 digits. N=450) gives the mean activity of the feature map $\bar{r}^l(\theta)$:

$$\bar{r}^l(\theta) = \frac{1}{N} \sum_{n=1}^{N} r^l(\theta, n) \tag{5}$$

Tuning values are defined for each feature map and digit, $d$ as:

$$f_d^l(\theta) = \frac{\frac{1}{N_d} \sum_{n \in d} r^l(\theta, n) - \bar{r}^l(\theta)}{\sqrt{\frac{1}{N} \sum_{i=1}^{N} (r^l(\theta, n) - \bar{r}^l(\theta))^2}} \tag{6}$$

When attention is applied to a particular digit, excitatory neurons that prefer that digit are given additional input. Specifically, the cells in feature maps whose tuning value for the attended digit are above the median tuning value for that digit are given attentional inputs. The attentional input to each feature map is proportional to how much above the median its tuning value is:

$$a_d^l(\theta) = \beta(f_d^l(\theta) - median(\mathbf{f}_d^l)) \tag{7}$$

Note, in this model the attentional input to the excitatory cell is fully specified by the above equation (that is, this value is not multiplied by the shape of the feedforward input).

40

We define digit preference on the feature map level (rather than for individual neurons) because feature attention is known to be a spatially-global phenomenon (that is, attention applied to a particular feature modulates neurons at all spatial locations, (Saenz et al., 2002)).

The accuracy on the same test set of overlaid images is again calculated for each digit, now in the presence of attention directed to the digit being detected. An additional parameter representing the overall strength of attention ($\beta$) is varied (.02, .04, or .06) and for each digit the best performing strength is used.

This attention was applied at each SSN layer individually as well as at both together. Here, the results of applying attention at the second SSN layer are reported as this elicited the best performance (a finding that is in line with those reported in (Lindsay and Miller, 2018; Lindsay, 2015), wherein attention at later layers better enhanced performance).

## 5. Acknowledgements

## References

D. Ress, B. T. Backus, D. J. Heeger, Activity in primary visual cortex predicts performance in a visual detection task, Nature neuroscience 3 (2000) 940.

S. Treue, J. H. Maunsell, Effects of attention on the processing of motion in macaque middle temporal and medial superior temporal visual cortical areas, Journal of Neuroscience 19 (1999) 7591–7602.

S. Treue, J. C. Martinez Trujillo, Feature-based attention influences motion processing gain in macaque visual cortex., Nature 399 (1999) 575–579.

M. R. Cohen, J. H. R. Maunsell, Attention improves performance primarily by reducing interneuronal correlations., Nat Neurosci 12 (2009) 1594–1600.

J. C. Martinez-Trujillo, S. Treue, Feature-based attention increases the selectivity of population responses in primate visual cortex, Current Biology 14 (2004) 744–751.

J. F. Mitchell, K. A. Sundberg, J. H. Reynolds, Differential attention-dependent response modulation across cell classes in macaque visual area v4., Neuron 55 (2007) 131–141.

J. H. Reynolds, D. J. Heeger, The normalization model of attention., Neuron 61 (2009) 168–185.

J. Lee, J. H. Maunsell, A normalization model of attentional modulation of single unit responses, PLoS ONE 4 (2009) e4651.

G. M. Ghose, Attentional modulation of visual responses by flexible input gain, J. Neurophysiol. 101 (2009) 2089–2106.

G. M. Boynton, A framework for describing the effects of attention on visual responses, Vision Res. 49 (2009) 1129–1143.

M. Carandini, D. J. Heeger, Normalization as a canonical neural computation, Nature Reviews Neuroscience 13 (2012) 51.

J. H. Reynolds, R. Desimone, Interacting roles of attention and visual salience in v4., Neuron 37 (2003) 853–863.

K. A. Sundberg, J. F. Mitchell, J. H. Reynolds, Spatial attention modulates center-surround interactions in macaque visual area v4., Neuron 61 (2009) 952–963.

Y. Ahmadian, D. B. Rubin, K. D. Miller, Analysis of the stabilized supralinear network, Neural computation 25 (2013) 1994–2037.

D. B. Rubin, S. D. Van Hooser, K. D. Miller, The stabilized supralinear network: a unifying circuit motif underlying multi-input integration in sensory cortex, Neuron 85 (2015) 402–417.

G. Hennequin, Y. Ahmadian, D. B. Rubin, M. Lengyel, K. D. Miller, The Dynamical Regime of Sensory Cortex: Stable Dynamics around a Single Stimulus-Tuned Attractor Account for Patterns of Noise Variability, Neuron 98 (2018) 846–860.

Y. Ahmadian, K. D. Miller, What is the dynamical regime of cerebral cortex?, arXiv preprint arXiv:1908.10101 (2019).

B. K. Murphy, K. D. Miller, Balanced amplification: a new mechanism of selective amplification of neural activity patterns., Neuron 61 (2009) 635–648.

G. W. Lindsay, K. D. Miller, How biological attention mechanisms improve task performance in a large-scale visual system model, eLife 7 (2018) e38105.

N. J. Priebe, D. Ferster, Mechanisms underlying cross-orientation supression in cat visual cortex, Nature Neurosci. 9 (2006) 552–561.

B. Li, J. K. Thompson, T. Duong, M. R. Peterson, R. D. Freeman, Origins of cross-orientation suppression in the visual cortex, J. Neurophysiol. 96 (2006) 1755–1764.

F. Sengpiel, V. Vorobyov, Intracortical origins of interocular suppression in the visual cortex, J. Neurosci. 25 (2005) 6394–6400.

H. Ozeki, I. M. Finn, E. S. Schaffer, K. D. Miller, D. Ferster, Inhibitory stabilization of the cortical network underlies visual surround suppression., Neuron 62 (2009) 578–592.

J. Lee, J. H. R. Maunsell, A normalization model of attentional modulation of single unit responses., PLoS One 4 (2009) e4651.

A. M. Ni, S. Ray, J. H. R. Maunsell, Tuned normalization explains the size of attention modulations., Neuron 73 (2012) 803–813.

G. Sclar, Expression of "retinal" contrast gain control by neurons of the cat's lateral geniculate nucleus., Exp. Brain Res. 66 (1987) 589–596.

G. Sclar, J. H. Maunsell, P. Lennie, Coding of image contrast in central visual pathways of the macaque monkey, Vision Res. 30 (1990) 1–10.

K. I. Naka, W. A. Rushton, S-potentials from luminosity units in the retina of fish (cyprinidae)., J Physiol 185 (1966) 587–599.

J. Martinez-Trujillo, S. Treue, Attentional modulation strength in cortical area mt depends on stimulus contrast., Neuron 35 (2002) 365–370.

M. Roberts, L. S. Delicato, J. Herrero, M. A. Gieselmann, A. Thiele, Attention alters spatial integration in macaque v1 in an eccentricity-dependent manner., Nat Neurosci 10 (2007) 1483–1491.

J. Lee, J. H. Maunsell, Attentional modulation of mt neurons with single or multiple stimuli in their receptive fields, Journal of Neuroscience 30 (2010) 3058–3066.

M. M. Churchland, B. M. Yu, J. P. Cunningham, L. P. Sugrue, M. R. Cohen, G. S. Corrado, W. T. Newsome, A. M. Clark, P. Hosseini, B. B. Scott, D. C. Bradley, M. A. Smith, A. Kohn, J. A. Movshon, K. M. Armstrong, T. Moore, S. W. Chang, L. H. Snyder, S. G. Lisberger, N. J. Priebe, I. M. Finn, D. Ferster, S. I. Ryu, G. Santhanam, M. Sahani, K. V. Shenoy, Stimulus onset quenches neural variability: a widespread cortical phenomenon., Nat Neurosci 13 (2010) 369–378.

L. Busse, A. R. Wade, M. Carandini, Representation of concurrent stimuli by population activity in visual cortex., Neuron 64 (2009) 931–942.

B. B. Averbeck, P. E. Latham, A. Pouget, Neural correlations, population coding and computation, Nature reviews neuroscience 7 (2006) 358.

D. A. Ruff, M. R. Cohen, Attention can either increase or decrease spike count correlations in visual cortex, Nature neuroscience 17 (2014) 1591.

S. Zhang, M. Xu, T. Kamigaki, J. P. Hoang Do, W. C. Chang, S. Jenvay, K. Miyamichi, L. Luo, Y. Dan, Selective attention. Long-range and local circuits for top-down modulation of visual cortex processing, Science 345 (2014) 660–665.

Y. Fu, J. M. Tucciarone, J. S. Espinosa, N. Sheng, D. P. Darcy, R. A. Nicoll, Z. J. Huang, M. P. Stryker, A cortical circuit for gain control by behavioral state, Cell 156 (2014) 1139–1152.

X. Jiang, G. Wang, A. J. Lee, R. L. Stornetta, J. J. Zhu, The organization of two new cortical interneuronal circuits., Nat Neurosci 16 (2013) 210–218.

M. E. Larkum, The yin and yang of cortical layer 1., Nat Neurosci 16 (2013) 114–115.

44

A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, pp. 1097–1105.

G. W. Lindsay, Feature-based attention in convolutional neural networks, arXiv preprint arXiv:1511.06408 (2015).

N. P. Bichot, M. T. Heard, E. M. DeGennaro, R. Desimone, A source for feature-based attention in the prefrontal cortex, Neuron 88 (2015) 832–844.

S. Paneri, G. G. Gregoriou, Top-down control of visual attention by the prefrontal cortex. functional specialization and long-range interactions, Frontiers in neuroscience 11 (2017) 545.

T. Miconi, R. VanRullen, A feedback model of attention explains the diverse effects of attention on neural firing rates and receptive field structure, PLoS computational biology 12 (2016) e1004770.

D. L. Yamins, J. J. DiCarlo, Using goal-driven deep learning models to understand sensory cortex, Nature neuroscience 19 (2016) 356.

A. J. Kell, J. H. McDermott, Deep neural network models of sensory systems: windows onto the role of task constraints, Current opinion in neurobiology 55 (2019) 121–132.

M. Kaschube, M. Schnabel, S. Lowel, D. M. Coppola, L. E. White, F. Wolf, Universality in the evolution of orientation columns in the visual cortex., Science 330 (2010) 1113–1116.

K. Albus, A quantitative study of the projection area of the central and the paracentral visual field in area 17 of the cat. i. the precision of the topography., Exp Brain Res 24 (1975) 159–179.

B. Haider, M. R. Krause, A. Duque, Y. Yu, J. Touryan, J. A. Mazer, D. A. McCormick, Synaptic and network mechanisms of sparse and reliable visual cortical activity during nonclassical receptive field stimulation., Neuron 65 (2010) 107–121.

G. Turrigiano, Too many cooks? intrinsic and synaptic homeostatic mechanisms in cortical circuit refinement., Annu Rev Neurosci 34 (2011) 89–103.

1116 M. Saenz, G. T. Buracas, G. M. Boynton, Global effects of feature-based
1117 attention in human visual cortex, Nature neuroscience 5 (2002) 631.
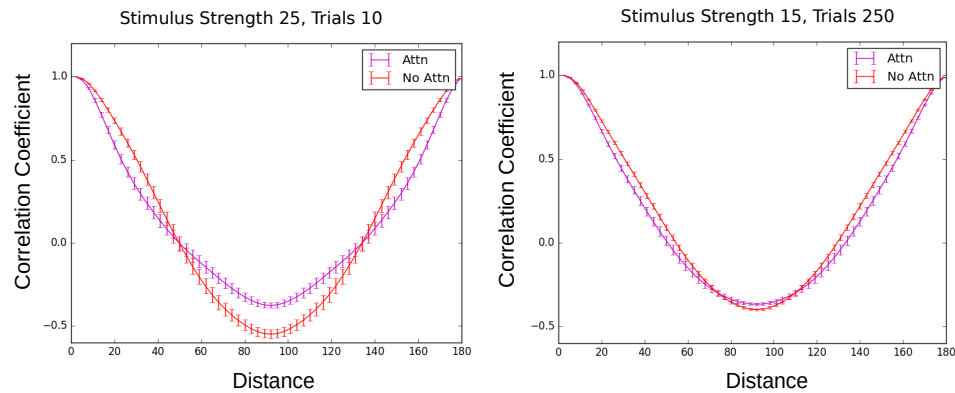
## 1118 Appendix A. Supplementary Figures



Figure A.17:

**Attention can increase correlations.** Example runs of the model used to make Figure 14 that result in attention increasing correlations for distant pairs. The strength of the stimulus and number of trials used for each condition is given at the top for each (in Figure 14, strength was 25 and 500 trials were used). Errorbars are SEM.
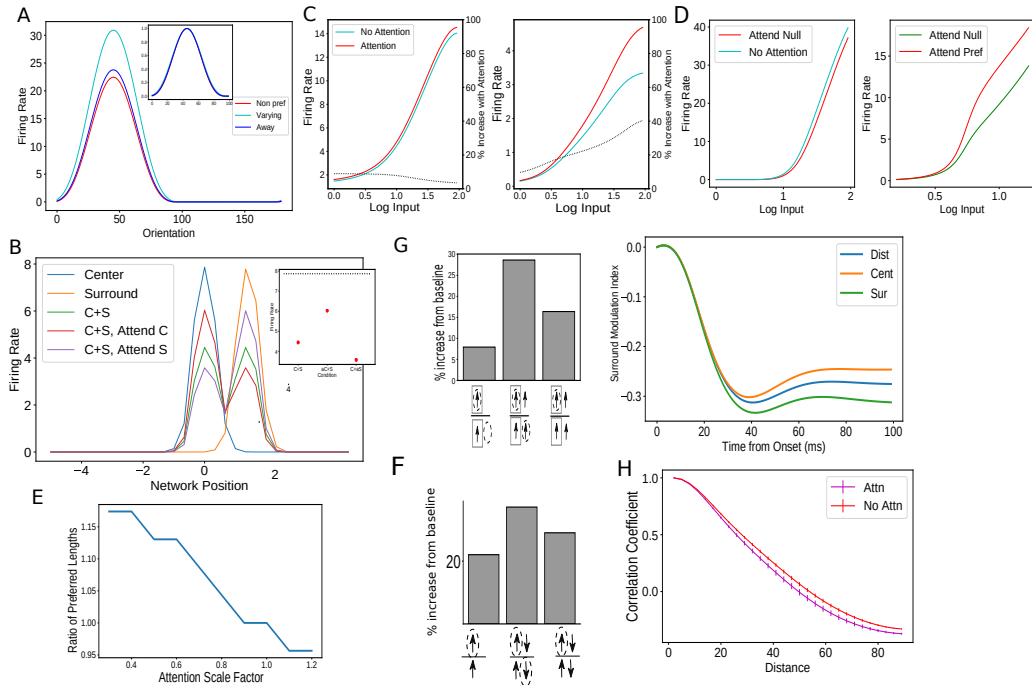
46

Figure A.18:

**Findings that qualitatively replicated with attention modeled as inhibitory input to inhibitory cells** A. Replication of Figure 4. B. Replication of Figure 6. C. Replication of Figure 8. D. Replication of Figure 9. E. Replication of Figure 10. F. Replication of Figure 11. G. Replication of Figure 12. H. Replication of Figure 14.
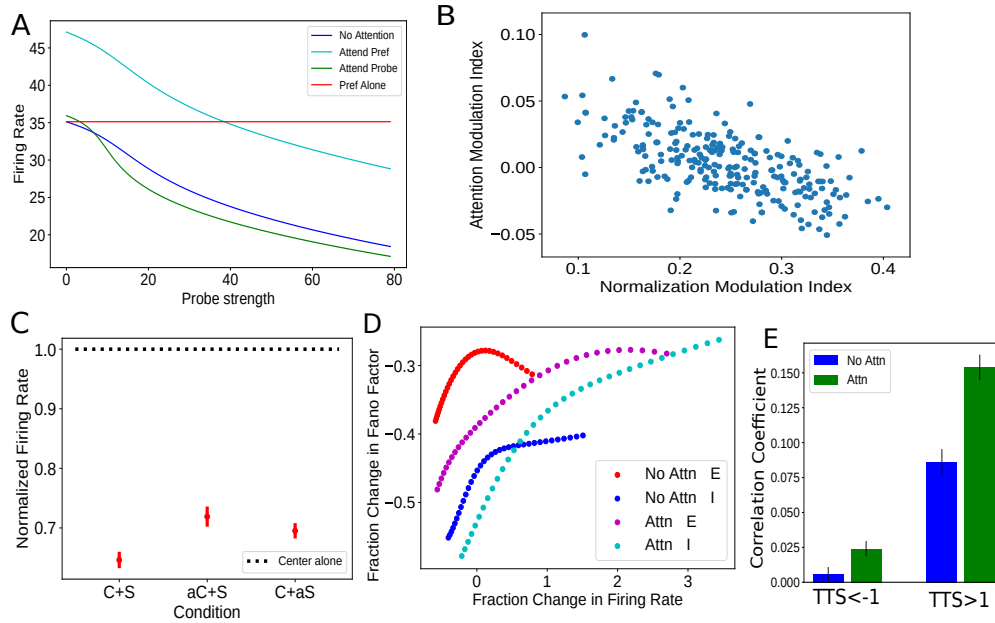
Figure A.19:

**Findings not qualitatively replicated with attention modeled as inhibitory input to inhibitory cells** A. Figure 3. Here much of the results are replicated however at low probe strengths attending the probe can increase firing rates compared to no attention. B. Figure 5. Here the relationship between normalization and attention is negative. C. Figure 7. Here the attend-surround condition is too similar to the attend-center one. D. Figure 13. Here for a range of firing rate changes, inhibitory cells have their Fano Factor increased with attention (though it should be noted this result happens occasionally when modeling attention as excitation to excitatory cells, for example, when the number of trials is lower). E. Figure 15. Here cell pairs with TTS¿1 also show an increase in correlation with attention.
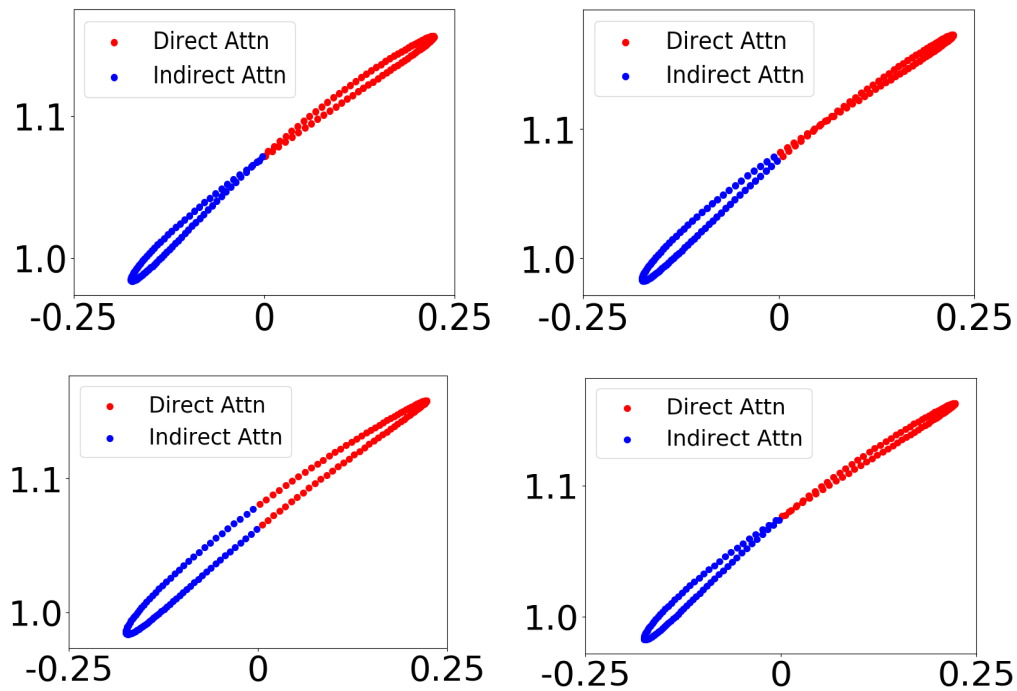
Figure A.20:

**Impact of feature attention at different spatial locations in layer 2 of the SSN-CNN** Ratio of attended to non-attended firing rates for cells in a ring network as a function of tuning value as in Figure 16E, but for different nearby spatial locations.