

Feedforward inhibition allows input summation to vary in recurrent cortical networks

Mark H. Histed

National Institute of Mental Health, National Institutes of Health, 35 Convent Dr. 35/3A-203, Bethesda, MD 20892. mark.histed@nih.gov

Brain computations depend on how neurons transform inputs to spike outputs. Because sensory stimuli change activity in many interconnected brain areas, it has been challenging to control neuronal inputs to measure input-output transformations *in vivo*. To overcome that difficulty, here we paired optogenetic stimuli with a constant sensory stimulus and measured spiking of visual cortical neurons in awake mice. We found that neurons' average responses were surprisingly linear. We then used a recurrent cortical network model to determine if these data and past observations of sublinearity could be described by a common circuit architecture. The model showed the input-output transformation could be changed from linear to sublinear with moderate (~20%) strengthening of connections between inhibitory neurons, but this change depends on the presence of feedforward inhibition. Thus, feedforward inhibition, a common feature of cortical circuitry, enables networks to flexibly change their spiking responses via changes in recurrent connectivity.

Introduction

Neurons in the cerebral cortex receive thousands of synaptic inputs and transform those inputs into spike outputs. Input-output transformations can be characterized in single cells (measuring firing rate while injecting current, producing a f-I curve, (Connors et al., 1982; Destexhe and Paré, 1999; Koike et al., 1970)), but network effects can dramatically alter input-output transformations *in vivo*. For example, ongoing network activity can create supralinearities in neurons' input-output functions (Priebe and Ferster, 2008), strong network connectivity can create entirely linear input-output functions (Brunel, 2000; van Vreeswijk and Sompolinsky, 1996), and recurrent connections can amplify inhibition to produce sublinearity (Ahmadian et al., 2013).

In this work, we examine input-output transformations *in vivo* by first measuring spiking responses to combinations of visual and optogenetic input in the mouse visual cortex (V1). Then, to shed light on the network and circuit mechanisms of input-output transformations, we use a spiking recurrent network model. The experimental data shows a strikingly linear input-output transformation in mouse V1, which stands in contrast to sublinearity seen in monkey V1 (Nassi et al., 2015). The model shows that the cortical network can achieve both kinds of transformations with only moderate changes in local recurrent synaptic strengths. The model makes a further prediction that feedforward inhibition – input that synapses not just on excitatory but also on inhibitory neurons – allows the cortex to support both kinds of transformations.

It is difficult to fully characterize input-output transformations using sensory stimuli alone, because sensory stimuli are processed by many brain regions each of which may provide input to a cortical area under study. Combinations of sensory stimuli have, however, found that a wide range of transformations are possible, often finding evidence for normalization, a form of sublinear summation (Carandini and Heeger, 2012). A few recent studies have used direct optogenetic input to study input-output transformations, and studies in different species have observed both normalization (Nassi et al., 2015; Sato et al., 2014) and more linear summation (Huang et al., 2014).

Models and theoretical approaches complement experimental studies of input-output transformations, because is difficult to control connectivity in an *in vivo* cortical network experimentally. Rate-based models (Ahmadian et al., 2013; Rubin et al., 2015) have characterized the range of behaviors cortical networks can support. But not all of the effects seen in rate-based models may occur in biological networks, as spiking neurons have biophysical properties that can impact input-output transformations such as refractory periods and nonlinearities due to spike threshold. Analysis of networks of spiking neurons is most advanced for models that approximate neuronal inputs as currents and not conductances (e.g. Brunel, 2000), but input-output relationships can be modified by the changes in effective synaptic strength and Vm variability (Richardson, 2004, 2007) that occur in realistic conductance-based neurons. Therefore, we use numerical simulations of models of conductance-based spiking neurons to determine which connectivity properties might create the input-output transformations seen in our data and in past data.

Below, we first describe the experimental results from mouse visual cortex (Fig. 1), showing near-linear responses across a wide range of firing rates and visual contrast. We then describe results from the model, showing that feedforward inhibition can produce sublinearity (Fig. 2), and that with feedforward inhibition, local connectivity can allow networks to be either linear or sublinear (Figs. 3-4). Finally, we construct a model network to fit our optogenetic-visual stimulation data (Fig. 5), which predicts that a canonical cortical circuit, with feedforward inhibition and varying only local connectivity, can resolve the differences in input-output transformations seen in V1 of different species.

Results

Experimental measurements in mouse V1 show linear summation

We combined visual and optogenetic input (Fig. 1A-B) by expressing ChR2 in V1 excitatory neurons using a transgenic mouse line and a Cre-dependent virus, and we used blue light pulses several seconds in duration (4-6 sec) to shift network firing rates to a new baseline. We delivered the same visual stimulus repeatedly, with and without ChR2 stimulation. We kept animals alert by giving them drops of fluid approximately once a minute, and we measured neurons' spiking via extracellular recording with multi-site probes.

When we presented the same visual stimulus with and without optogenetic stimulation, we found that V1 neurons' responses scaled nearly linearly (Fig. 1C) – that is, the same size response was produced even as the optogenetic stimulus changed the baseline firing rate. Even for relatively large optogenetic baseline shifts (~10 spk/s, roughly the same magnitude as the average visual response), the visual response was similar with and without ChR2 stimulation. We saw this linear response across a range of intensities of the visual stimulus (contrast range: 8%-90%, Fig. 1D), and we saw linear responses both when averaging single units (N=50) and multi-units (N=239). Responses became slightly sublinear in cells with the largest baseline shifts (Fig. 1E), but responses were on average within a few percent of linear (for maximum contrast, as in Fig. 1D: response change for single units -4.8%, for multi-units -4.1%).

While average neuronal responses were nearly linear, individual recorded units were often either supra- or sub-linear. For example, in Fig. 1E, many points lie above and below the horizontal line that shows a perfectly linear response. (Several example units that produce either sub- or supralinear responses are shown in Supp. Fig. 1.) With the 90% contrast visual stimulus, 34% of single units are significantly non-linear (17/50, $p < 0.01$, KS test; Supp. Fig. 1C), and 28% of multi-units are significantly non-linear (67/239, $p < 0.01$, KS

test). Such heterogeneity in responses could arise because each neuron has slightly different local connectivity. Heterogeneity due to local recurrent connections would suggest the population average linear response is a network effect, arising from connections between excitatory and inhibitory neurons that cause them to dynamically respond to each others' activity (van Vreeswijk and Sompolinsky, 1996). Below, we test how connectivity might lead to the observed responses, using a spiking network model.

Other experimental work finds sublinear summation in macaque visual cortex

In contrast to this average linear scaling in mouse primary visual cortex, recent work in the monkey primary visual cortex (Nassi et al., 2015) found neural responses that were at times highly sublinear, and sublinear on average. The experimental approach used by Nassi et al. does not seem to differ in important ways from our approach -- they expressed ChR2 primarily in excitatory neurons using the CaMKII-alpha promoter, stimulated an area of the cortex a few hundred microns in diameter, and they paired ChR2 and visual stimulation. Because the different results may stem from differences in cortical architecture across species, rather than differences in experimental methods, we sought to determine whether there were features of local cortical circuits that could change response scaling from linear to sublinear.

Model network simulations identify circuit properties controlling input summation

Since it is difficult to manipulate neural connectivity *in vivo*, we used numerical simulations of conductance-based model neurons to understand how network connectivity might change response scaling. We constructed networks of 10,000 conductance-based leaky integrate-and-fire neurons, 8,000 excitatory (E) and 2,000 inhibitory (I). We chose realistic parameters for the model neurons, including sparse connectivity (initially 2%), and chose moderate synaptic strengths such that a few tens of EPSPs were required to push a neuron over threshold. (We explore a range of values of sparsity and synaptic strength below.) These sparse, randomly connected networks produce irregular and asynchronous spontaneous activity (Fig. 2A) similar to that observed experimentally (Destexhe et al., 2003; Steriade et al., 2001) and show stable responses to external inputs (Vogels and Abbott, 2005). For all simulations, we set the spontaneous average rate of the network to 5 spk/s. There are a variety of single-cell properties that could set neurons' spontaneous rate, but we changed the spontaneous rate by supplying a small, constant amount of excitatory input (that does not vary with network activity or input) to either excitatory or inhibitory neurons (see Methods).

To determine how different sorts of feedforward inputs affect neurons' responses, we simulated external inputs to E and I cells by two input groups of Poisson spike trains whose rates could be varied independently. As expected, when we varied the external input rates, increasing input to E cells (x-axis) monotonically increased the average network response (Fig. 2B, contour lines; average of all excitatory cells in the network, a measure similar to that obtained by multi-electrode recordings) and increasing input to I cells (y-axis) monotonically decreased the average network response. However, we could hold the average response constant by adjusting the two feedforward inputs. When the average response was constant (along contour lines in Fig. 2B), we still observed changes in response scaling, and those changes depended on the amount of I input.

To assess response scaling, we began with a combination of E and I input that produced a 15 spk/s response (chosen because experimentally, we measured an average response that peaked near 15 spk/s, Fig. 1C,D). Then, we multiplied both input rates by a single constant and measured the size of the response to the scaled input. We found that when feedforward I input is small, responses are near-linear (Fig. 2C). This is not surprising, as previous theoretical work using strong local synaptic coupling in models with binary (van Vreeswijk and Sompolinsky, 1996) or current-based neurons (Brunel, 2000) showed that networks can produce linear responses even though individual neurons in cortical networks are nonlinear (Priebe and Ferster, 2008). However, these models did not characterize the effects of varying feedforward E and I input separately, and so we varied feedforward I input in the conductance-based model. Indeed, when feedforward I input was varied, we observed deviations from linearity. Even though the spontaneous spike rate and the spike rate response to a single stimulus alone were both held constant with and without feedforward inhibition, increasing stimulus strength showed more sublinear response scaling when feedforward inhibition was present.

Local connectivity changes summation only in the presence of feedforward inhibition

While adding feedforward inhibition induced some sublinearity, we wished to know if more dramatic nonlinearities were possible. Therefore, we next changed local recurrent connectivity between and amongst E and I populations, and measured how those connectivity changes affected responses scaling (Fig. 3). In Fig. 3, we show the effects of varying two local connections (first, strength of synapses from E to I, and second, strength of synapses from I to I) to illustrate the range of effects we observed. (Supp. Fig. 2 shows the effects of varying all pairwise combinations of E to I connectivity, as well as feedforward E and I input strength.) To implement varying connectivity in the model, we added additional connections between two neuronal populations (e.g, E to I, or I to I) with the same sparsity as the network. We then varied the strength of those additional connections and measured effects on response scaling.

With only feedforward input to E cells (Fig. 3A,C,E), we found that changing network connections did not dramatically affect response scaling. Changing the connectivity could change the gain of the network (the size of the response to a constant input, Fig. 3A, contour lines), but response scaling was nearly linear (Fig. 3A, plot is yellow throughout; Fig. 3C-D: black lines lie close to horizontal dotted line). At high firing rates, we consistently saw moderate increases in sublinearity, which seems likely to be due to effects of the absolute refractory period (3 ms) and thus, rates above 50 spk/s are indicated by light gray lines (Fig. 3CD). But the linear scaling we had observed in the model when delivering input to E cells only was robust to changes in local connectivity. In sum, without feedforward inhibition, scaling was approximately linear, and local connectivity changes had little effect.

Even though near-linear scaling was consistently seen when feedforward input arrived to E cells, when feedforward input arrived to both E and I cells, responses could be either linear or dramatically sublinear. When we increased local I to I connection strength (Fig 3B, y-axis), sublinearity was observed (Fig. 3D; plot parameters correspond to pink asterisk in Fig. 3B, in blue region of plot). But increased E to I connection strength (Fig. 3B, x-axis) led to more-linear scaling (Fig. 3E; plot parameters correspond to pink ‘|’ symbol in Fig. 3B). The

sublinear scaling produced by stronger I to I connectivity was dramatic. As with all the timecourse plots (Fig. 3C-F), we chose input strength so the first firing rate response was 15 spk/s, but when I to I connectivity was increased, subsequent firing rate responses fell as low as 1 spk/s (Fig. 3D). The mechanism by which increased I to I coupling produces increased sublinearity is not yet understood. Such unintuitive changes can arise from network-level effects, similar to the way E-I tracking may cause inhibitory neurons to actually decrease their activity when inhibitory neurons are excited by stimulation (Ahmadian et al., 2013), or may arise from cell-autonomous changes in conductance that leads to shunting in individual cells (Chance et al., 2002; Richardson, 2004). Further theoretical work will be required to understand why increased I-I coupling leads to increased sublinearity in spiking networks. However, it is likely that cortical inhibitory neurons *in vivo* do have the capacity to adjust their local connectivity, as inhibitory cells modify their dendritic structure over time (Chen et al., 2011). In sum, the numerical simulations show that local connectivity changes can dramatically affect response scaling, but only in the presence of feedforward I input.

Connectivity effects on summation do not depend on connection sparsity or strength

We next examined whether synaptic strength and connection sparsity can change the role of feedforward inhibition in response scaling. We expected that varying the total recurrent input that neurons receive would change non-linearity of responses, as predicted by theory (Ahmadian et al.; van Vreeswijk and Sompolinsky, 1996), as long as the network remained stable. Therefore, we varied total input in two ways, by varying connection sparsity and by varying synaptic strength (Fig. 4). Experimental estimates of local connection sparsity can range as high as 10-20% (*i.e.* each neuron connects to 10-20% of nearby neurons, Braitenberg and Schüz, 2001; Lefort et al., 2009). But the effective sparsity of connections might be lower, as connection probability in cortical networks is known to fall off with distance, so averaging connection probabilities across the network can give lower values than measured for nearby pairs. To examine the effects of changing connection probability, we varied sparsity between 2-20% and found that in all cases, adding feedforward inhibitory drive allowed more sublinear responses (Fig. 4; green lines always lie below blue lines in Fig 4A). We observed more linear scaling when we increased the strength of all synapses together (Fig. 3), and a bigger range of possible scaling (from supralinear to sublinear) when we decreased synaptic strength. These results show that, in networks that use a range of connection strength and sparsity, feedforward inhibition enables local E and I connectivity to have similar effects on response scaling, though the networks become more linear as connectivity strength increases.

Next, we asked whether a model fit to our ChR2 experiment shows similar scaling dependence on feedforward inhibition. Up to this point, we have examined the behavior of simulated networks only by scaling a feedforward input (Figs. 2-4). We have implemented this feedforward input to simulate the way input spikes change conductance in neurons, by modulating the firing rate of a (Poisson) stochastic point process. Using these input spike trains, the sum of feedforward synaptic inputs in a given network neuron has substantial fluctuations about its mean. In contrast, experimental ChR2 stimulation activates many channels, and produces conductance changes with much smaller fluctuation about the mean. Thus, it might be possible that the scaling behavior we studied experimentally, with

ChR2 combined with visual stimuli, would differ from the combinations of feedforward input we simulated in Figs. 2-4. To determine if there was a difference, we simulated ChR2 input by changing conductance and combined this with feedforward input (Fig. 5), and found that combinations of ChR2 and visual inputs produced qualitatively similar effects to the effects we had previously seen. Combinations of simulated ChR2 and visual input (Fig. 5A) showed slightly increased sublinearity when compared to a single scaled visual input (Figs. 2-4). (We also saw some slight sublinearity in our measurements of responses to combined ChR2 and visual input in mouse V1, Fig. 2.) However, as with simulated visual input (Figs. 2-4), we found that with paired conductance (ChR2) and spiking (visual) inputs, more sublinearity is possible when the feedforward input combines inhibitory and excitatory targets than when feedforward input targets only excitatory neurons (Fig. 5B-C). And, in the presence of feedforward inhibition, moderate changes in network connectivity can modify scaling behavior (Fig. 5D). In sum, in the models that use simulated visual (spiking) inputs of varying rate (Figs. 2-4), and the models that use combined visual and ChR2 (conductance) inputs (Fig. 5), the role of feedforward inhibition and I-I connectivity in response scaling is similar.

Summation in our data and past data can be explained by a model with feedforward inhibition

Finally, we constructed a model with combined visual (spiking) and ChR2 (conductance) inputs, and fit evoked rates to our data. We then asked what combinations of connectivity and feedforward input could describe both our data and past measurements. Our data (Fig. 5E) was well-matched by the simulations that showed small sublinearity (Fig. 5A-D). The data was similar to two different sets of network simulation parameters. Networks with only feedforward excitation showed responses that paralleled the data, but we also saw effects that paralleled the data in networks with both feedforward excitation and inhibition, for particular values of local connectivity. Since feedforward inhibition is a common feature of cortical networks in many species (Douglas and Martin, 2004), a model using feedforward inhibition seems a good choice to describe experimentally measured response scaling. With feedforward inhibition, changes in local (e.g. I-I) connectivity can change response scaling from linear to sublinear. While other network architectures might also give sublinear scaling, these simulations show that a wide regime of cortical scaling behavior, from linear (as seen here in mouse V1 and also in the tree shrew (Huang et al., 2014)), to strongly sublinear (as seen in primate V1, Nassi et al., 2015), can be achieved by a model with feedforward inhibition. In sum, the simulations show that model with fixed input connectivity, with feedforward inhibition, can describe both our data (Fig. 5E) and past observations.

Discussion

We have experimentally measured the average firing rate of a group of cortical neurons while presenting the same visual stimulus repeatedly. We found that average response summation in mouse V1 is close to linear, even though individual cells can be nonlinear. Linear summation holds even for substantial shifts in firing rate (ChR2-induced firing rate

changes of 10-15 spk/s, approximately the same size as the maximum visual response, Fig. 1). Using a numerical model of conductance-based spiking neurons, we find that response scaling is affected dramatically by synaptic connectivity. Moderate changes in synaptic coupling (~20%) between inhibitory cells can change response scaling from linear to sublinear (Figs. 3-5). Further, the change in inhibitory-to-inhibitory connectivity that leads to sublinear summation only yields such sublinear summation in the presence of feedforward inhibition.

Several types of input-output transformations have been characterized in brain circuits. Neuronal responses to multiple sensory stimuli can often be governed by normalization, where adding an additional stimulus yields divisive reduction of the responses to a single stimulus. This form of sublinear summation has been observed in different visual cortical areas of several species (Carandini and Heeger, 2012). Linear summation, on the other hand, is also commonly seen at various stages of sensory systems (Carandini and Heeger, 2012) and both linear and sublinear responses may be useful at different levels. Linear summation may be more desirable when responses at different locations should receive equal weight, as when an organism must sensitively detect a distant predator, or when spikes that occur at different times should produce the same downstream effect. In fact, some high-performance computer vision systems use both linear and normalization steps in distinct layers or networks (Carandini and Heeger, 2012; Yamins and DiCarlo, 2016). Experimentally, normalization is usually measured in the context of sensory stimuli, not with direct cortical input, and thus normalization might partially depend on subcortical (e.g. thalamic gain control, Bonin et al., 2006) or feedback effects. Advances in optical stimulation promise to allow fuller characterization of input-output transformations in brains *in vivo*.

Many neurons in the cortex change their firing rate in response to even small sensory stimuli (Bonin et al., 2011; Van Essen et al., 1984). Anatomically, sensory input that arrives to multiple cells is common, as in the case of divergent feedforward thalamic input to the cortex (Reid, 2001). Single axons from the thalamus often ramify across several hundred microns of the cortex (Braitenberg and Schüz, 2001; Garraghty and Sur, 1990), and thalamic axons projecting to the visual cortex can make synapses on dozens of excitatory cortical cells (Freund et al., 1989). Therefore, we delivered optogenetic input to multiple neurons simultaneously (using a blue light spot a few hundred μm in diameter, comparable to the region of mouse V1 activated by the small Gabor visual stimulus we used).

Optogenetic stimuli may lead to firing rate changes in other parts of the brain besides the area stimulated. But perhaps because the majority of synapses made by cortical neurons are within the same cortical area, local intracortical effects for optogenetic stimuli like these have been observed to be larger than remote effects on the visual thalamus (Li et al., 2013; Olsen et al., 2012). This is true even though the visual thalamus (dorsal lateral geniculate) receives a large proportion of all projections out of V1 (Reid, 2001). Thus, the neurons best suited to act as the recurrent population in the model may be other V1 neurons, and perhaps even neurons within a few hundred microns of the neurons receiving input, where the probability of recurrent connectivity is highest (Lefort et al., 2009). However, we do not rule out the possibility that other neurons in the brain contribute to the recurrent population.

Our results show that network mechanisms can contribute to response summation. The model neurons are leaky integrate-and-fire neurons, so individual neurons sum their subthreshold inputs entirely linearly, and the nonlinear spiking responses we characterize likely arise from how E and I neurons interact. We chose this model architecture because we judged it the simplest model that could capture both excitatory-inhibitory interactions and also single-cell nonlinearities due to refractory period, Vm fluctuations, spike threshold, and conductance changes (Chance et al., 2002; Richardson, 2004). There are, however, other single-cell mechanisms, such as short-term synaptic plasticity or dendritic nonlinearity (Häusser et al., 2000; Silver, 2010) that might additionally contribute to even more nonlinear summation, both below threshold and in spike responses. On the other hand, dendritic nonlinearities might also have roles that do not affect scaling, as for example nonlinearities can be used to amplify distant input synapses so that different synapses produce equal responses at the soma (Katz et al., 2009). These additional mechanisms could amplify or otherwise modify the network effects we have observed, and if those mechanisms are used they could also vary across species.

We adjusted synaptic coupling between (E and/or I) populations by changing the strength of a set of fixed connections between the desired populations. Because in sparse networks like this, neurons share only a small fraction of their input, we expected increases in synaptic strength to achieve the same qualitative result as adding new synapses, even if the two types of changes may not have exactly proportional effects on the behavior of the network. Fig. 4 shows that feedforward inhibition allows more sublinearity across changes in both synaptic strength and synapse number.

The linear responses we observed in mouse primary visual cortex are similar to those seen in tree shrew visual cortex (Huang et al., 2014), but are different than the sublinear responses seen in macaque visual cortex (Nassi et al., 2015). Our simulations show that a broadly similar cortical architecture can support both kinds of scaling of feedforward input, subject to moderate adjustments in local connectivity. The linear responses we saw in the mouse differ from those of Sato et al. (Sato et al., 2014), who also delivered combinations of optogenetic and visual input to mouse V1 neurons and found sublinearity under certain conditions. However, Sato et al. used an experimental approach different than the other three studies (macaque, tree shrew and our study in mouse), in which they optogenetically elicited antidromic input spikes by stimulating the contralateral hemisphere from which they were recording. Comparing these two types of input may shed additional light on how cortical circuits transform inputs to outputs.

Feedforward inhibition is included in the canonical cortical microcircuit framework (Douglas and Martin, 2004) because it is a stereotypical feature of many cortical areas. In sensory cortical areas, including the visual cortex, it has been observed that input thalamic neurons make synapses both onto excitatory principal cells and onto inhibitory basket cells. Such feedforward inhibitory connectivity has been observed both with anatomical and physiological methods (Isaacson and Scanziani, 2011). Since inhibitory basket cells project strongly back to excitatory cells, inhibitory changes due to thalamic input arrive to principal cells a few milliseconds after the first excitatory changes. This delay of a few milliseconds between the arrival of excitation and inhibition can be used to align spike outputs of cortical neurons (Cruikshank et al., 2007; Gabernet et al., 2005; Swadlow, 2003; Tiesinga et al.,

2008). Beyond shaping the timing of spike responses, however, it has been previously noted that feedforward inhibition might also be used to control the magnitude of spiking responses to thalamic input. Douglas et al. (Douglas et al., 1995) proposed that spike responses can be shaped by preferential amplification of either excitation or inhibition in cortical recurrent networks, where amplification might arise by connections within populations of excitatory or inhibitory neurons. Ahmadian and Miller (2013) later showed that rate-based networks with an excitatory and inhibitory term that are stable (so that the network does not e.g. diverge and become epileptic) have regimes of both linearity and sublinearity, although it is not yet clear which of these regimes spiking networks operate in, and which cellular or synaptic parameters affect summation. In Ahmadian and Miller's model, individual cells can be supralinear (Priebe and Ferster, 2008), but when external drive arrives to multiple cells, supralinearity is seen only when recurrent connections are weak and thus excitation and inhibition are not strongly coupled. This may explain why we saw supralinear responses only in the model network with the weakest synaptic connectivity (Fig. 4).

Substantial recurrent intracortical responses are elicited by sensory input, with approximately $2/3^{\text{rd}}$ of synaptic input after a sensory stimulus arising from recurrent synapses (Li et al., 2013; Lien and Scanziani, 2013). With strong recurrent connectivity, previous modeling results (Renart et al., 2010; van Vreeswijk and Sompolinsky, 1996) predict that excitatory and inhibitory populations are forced by the strong coupling to track each others' activity closely, resulting in linear responses. In accord with this prediction about strongly-coupled networks, we observed increasing linearity when we increased synaptic strength (Fig. 4) as long as the network remained stable. However, for very strong recurrent connectivity, feedforward connectivity must also be very strong to drive any response (Ahmadian et al., 2013)(see also our Fig. 4), which appears non-physiological (Li et al., 2013; Lien and Scanziani, 2013). Most of our simulations (Figs. 2,3,5 and parts of Fig. 4) use synapses of moderate size (order 1mV, see Methods), requiring tens of PSPs to combine to produce a spike, as seen in cortical neurons (Barral and Reyes, 2016). Taken together, these observations imply the differences in scaling we observed occur in a range of moderate synaptic strengths: low enough to avoid obligate linearity, and high enough to allow recurrent connections to contribute substantially to network input-output functions.

We found that a network model can link local connectivity to network physiological responses in ways that might be difficult to predict without the model. It has been difficult to measure many of the synapses in a brain volume, but connectomic methods such as those using large-scale electron microscopy (Briggman et al., 2011; Lee et al., 2016) promise to make such comprehensive synaptic mapping possible even in column-sized volumes of the cortex. Combining approaches for controlling input with methods to measure connectivity will be useful to shed light on an important part of brain computation – the input-output transformations of populations of connected cells.

Methods

Neurophysiology

All experimental animal procedures were conducted in accordance with NIH standards and were approved by the IACUC at Harvard Medical School. Animal breeding and surgery were performed according to the methods described previously (Glickfeld et al., 2013; Histed and Maunsell, 2013).

Neurophysiological data from Emx1-Cre animals (N=4) were collected using the methods used in Glickfeld et al. (2013) for extracellular recordings. Briefly, animals kept on a monitored water schedule were given small drops of water (~1 μ L) every 60-120 s during recording to keep them awake and alert. The visual stimulus was presented for 200 ms with a blank period of 800 ms between presentations. Gabor patches had spatial frequency of 0.1 cycles/deg and sigma of 12.5 deg. Optogenetic light pulses were delivered on alternating sets of 10 stimulus presentations (light onset 500 ms before first stimulus, offset 500ms after end of last stimulus; total light pulse duration 10.2s). A 1 s delay was added after each set of 10 stimulus presentations. Extracellular probes were 32-site silicon electrodes (Neuronexus, Inc., probe model A4x8). Recording surfaces were treated with PEDOT to lower impedance and improve recording quality. On each recording day, electrodes were introduced through the dura and left stationary for approximately 1 hour before recording to give more stable recordings. As in (Histed and Maunsell, 2013) ChR2 was expressed in excitatory neurons using viral injections into the Emx1-Cre (Gorski et al., 2002), (Stock #5628, Jackson Laboratory, Bar Harbor, ME USA) line. Virus (0.25-1.0 μ L) was injected into a cortical site whose retinotopic location was identified by imaging autofluorescence responses to small visual stimuli. Light powers used for optogenetic stimulation were 500 μ W/mm² on the first recording session; in later sessions dural thickening was visible and changes in firing rate were smaller, so power was increased (maximum 3 mW/mm²) to give mean spontaneous rate increases of approximately ~5 spikes/s in that recording session. Optogenetic light spot diameter was 400-700 μ m (FWHM) as measured by imaging the delivered light on the cortical surface. Spike waveforms were sorted after the experiment using OfflineSorter (Plexon, Inc.). Single units were identified as waveform clusters that showed clear and stable separation from noise and other clusters, unimodal width distributions, and inter-spike interval histograms consistent with cortical neuron absolute and relative refractory periods. Multiunits were clusters that were distinct from noise but did not meet one or more of those criteria, and thus these multiunits likely group together a small number of single neuron waveforms.

Data analysis

Spike histograms were smoothed using piecewise splines (LOWESS smoothing). To compute the visual response for each neuron in Fig. 1D, we counted spikes over a 175 ms period beginning 25 ms after stimulus onset, with a matched baseline period 175 ms long ending at stimulus onset. The dataset includes data from 100 shank penetrations (~25 recording sessions with a 4-shank electrode). Because the inter-shank spacing was 200-400 μ m, our stimuli in fixed retinotopic locations could not activate neurons on all shanks. Therefore, we included only shanks in which an average visual response > 0.2 spikes/s was measured (38/100 shanks). This gave 417 single and multi-units. We examined only units that showed a visual stimulus response (N=289; mean stimulus response-mean spontaneous > 0.2) in the absence of ChR2 stimulation. Because ChR2 expression was highest at the site of viral injection and fell off with distance, we

took advantage of this variation to sort units into three groups based on the strength of local ChR2 activation (Fig. 1C). We found the average change in spontaneous rate induced by ChR2 stimulation for all units on a shank and rank-ordered the shanks. Dividing shanks into three groups based on small, medium, or large ChR2 effects yielded three nearly-equal sized groups of units receiving small, medium or large ChR2 activation. The group sizes differ by a few units because we sorted by shank, not by individual unit.

Conductance-based spiking network model

The cortical model is a recurrent network of conductance-based leaky integrate-and-fire neurons. Example Python code and a Jupyter notebook (<http://jupyter.org>) are provided at <https://github.com/histedlab/code-feedforward-inhibition-condLIF> that run the network simulation with all its inputs, replicating spike counts shown in Fig. 5C, bottom row. To recover the rest of the simulations in Fig. 2-5, this code can be run in parallel on a larger cluster.

Each model neuron is connected randomly to each other neuron with fixed probability (sparsity). For example, for a 10% sparsity network, each cell receives input from 10% of the excitatory cells and thus gets $0.1 \times 8000 = 800$ E inputs. Similarly, at 10% sparsity, each cell receives $0.1 \times 2000 = 200$ I inputs. As seen in the cortex, we chose the inhibitory synaptic strength to be larger than the excitatory synaptic strength to achieve rough balance, but we also varied both synaptic strengths and found that our conclusions are not affected by changes in E/I synaptic strength ratio. (Supp. Fig. 2; see also Fig. 4 for effects of changing together E and I recurrent synaptic weights by an order of magnitude). We refer to this baseline set of random, sparse connections as the balancing connections. For convenience, to change local connectivity, we change the strength of a second added set of connections with the same sparsity while keeping the strength of the balancing connections constant. For example, when I->I connectivity is varied in the 2% sparsity network (e.g. Fig. 3), each I cell receives an extra 40 synapses from other I cells, and the y-axis in Fig. 3AB shows the effects of varying the weight of those 40 synapses from zero to ~20% of the weight of the standard recurrent I->I synapses.

Each simulated neuron's membrane potential evolves according to the following equation:

$$\frac{dV_m}{dt} = -\frac{1}{\tau_m} \left[g_{leak}(V_m - E_{rest}) + g_{ChR2}(V_m - E_e) + g_e(V_m - E_e) + g_i(V_m - E_i) \right]$$

When the membrane potential V_m crosses a threshold (-50 mV), a spike is recorded and V_m is reset to E_{rest} (-60 mV) for the absolute refractory period (3 ms).

Beyond the recurrent inputs from other neurons in the network (described in the model architecture above), model neurons can receive two kinds of external inputs: external feedforward inputs simulating e.g. sensory input from thalamus, and external ChR2 inputs. Feedforward (sensory) inputs are simulated as Poisson spike trains whose rates are changed by stepping to a new value, with values chosen to approximate visually-evoked changes seen in the

data. ChR2 input is simulated by linearly ramping g_{ChR2} to a new value over 2 ms, a timescale consistent with ChR2 t_{on} (Nikolic et al., 2009), and g_{ChR2} amplitude is varied to reproduce experimental changes in firing rate (see below). Synaptic conductances g_e and g_i are incremented instantaneously by a constant excitatory or inhibitory synaptic weight when a spike is fired by a recurrent or feedforward input. The conductances decay with time constants $\tau_{ge} = 5\text{ ms}$ and $\tau_{gi} = 10\text{ ms}$, described by:

$$\frac{dg_e}{dt} = -\frac{g_e}{\tau_{ge}}$$

$$\frac{dg_i}{dt} = -\frac{g_i}{\tau_{gi}}$$

Other constants are: excitatory reversal $E_e = 0\text{ mV}$, inhibitory reversal $E_i = -80\text{ mV}$, membrane time constant $\tau_m = 20\text{ ms}$. Post-synaptic potential (PSP) amplitudes can vary with network activity and synaptic weight because the model neurons are conductance-based. As we varied sparsity in the network, the excitatory PSP amplitude (Fig. 3) varied over an approximately tenfold range (0.3-3.0 mV for sparsity 20% - 2%, if calculated assuming that the mean membrane potential of network neurons is -65mV.)

Network spontaneous firing rate

The sparse recurrent connections yield spontaneous activity in the network in the absence of external input (van Vreeswijk and Sompolinsky, 1998; Vogels and Abbott, 2005). To equate the spontaneous firing state of the network across different sparsity and synaptic strength, we adjust network spontaneous rate. We use an additional external Poisson excitatory input to either excitatory or inhibitory neurons to respectively raise or lower the spontaneous rate. The rate of this Poisson input is chosen via stepwise optimization to give a mean spontaneous rate across excitatory neurons of 5 spk/s. (In the 2% sparsity network, these added excitatory synapses account for only approximately 2% of the total mean conductance). For many networks, a local minimum of the parameter can be found repeatably, but for extreme values of sparsity and synaptic strength, the network is unstable and spontaneous rates are either sensitive to small perturbations or diverge. In these cases network response is not shown (e.g. gray regions, Fig. 4B-C).

Simulations were performed with the Brian simulator (Brette et al., 2007) on a multi-CPU cluster (the NIH HPC Biowulf cluster, <http://hpc.nih.gov>, or Orchestra, <http://rc.hms.harvard.edu>) with an integration time step of 50 μs .

Acknowledgements

This work was funded in part by the Intramural Research Program of the NIMH and by U01 NS090576 (BRAIN Initiative). We thank Nicolas Brunel for discussion and comments, Alex Handler and Oliver Mithoefer for technical assistance with neurophysiology, John Maunsell for support, and Bruno Averbeck, Barry Richmond, Alessandro Sanzeni, Lex Kravitz, and Carson Chow for comments on the manuscript.

Figures

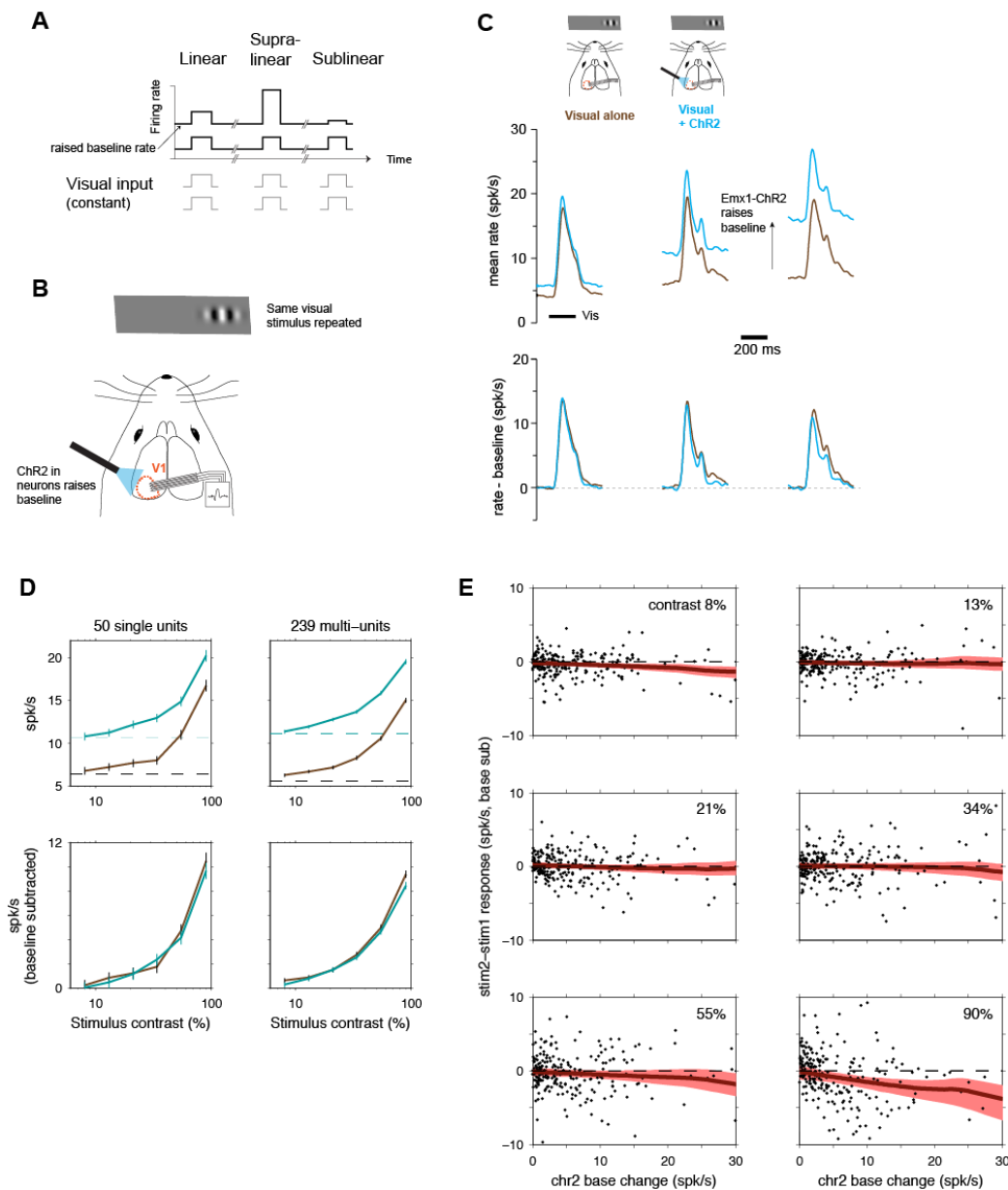


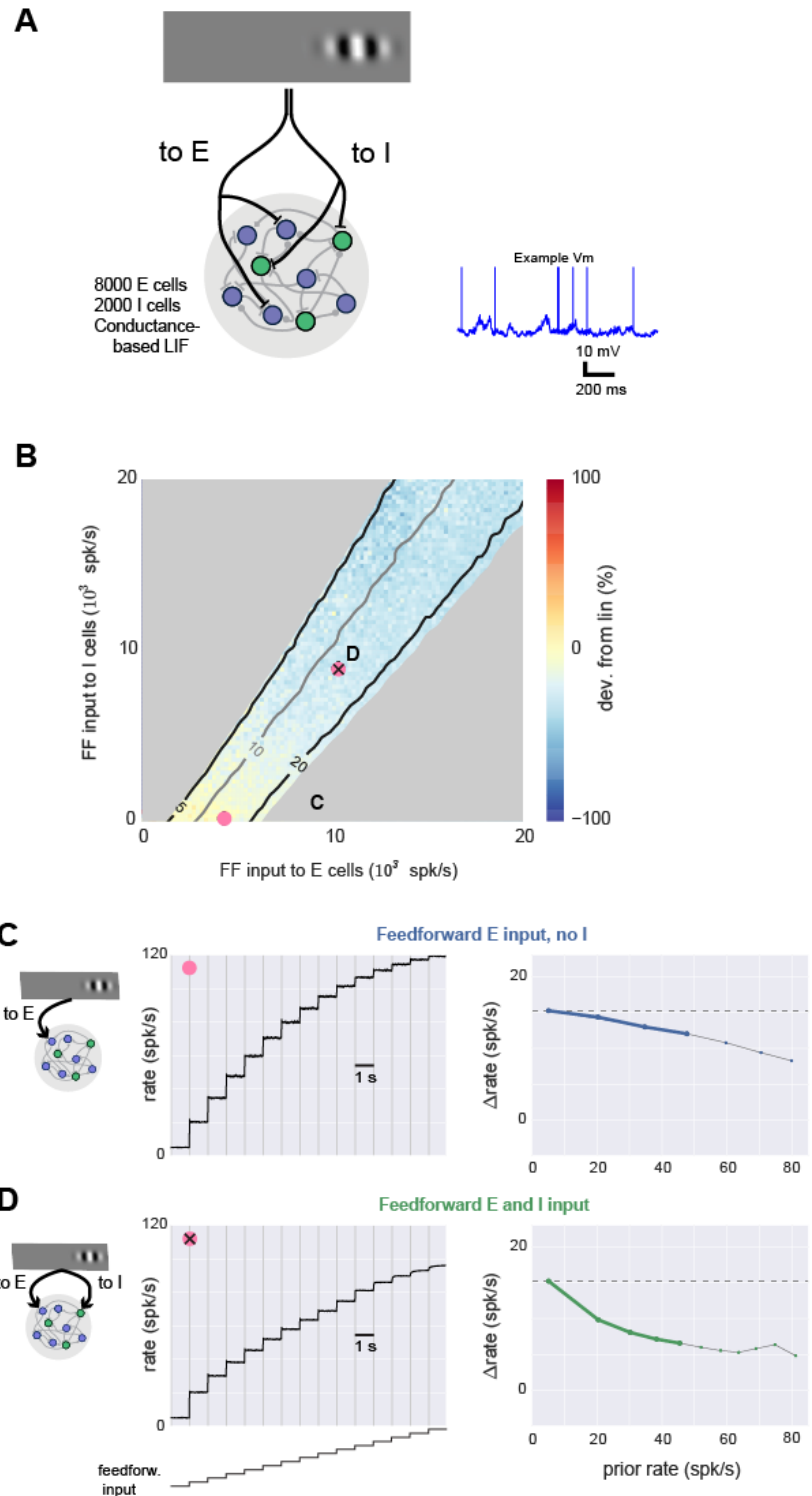
Figure 1: Near-linear scaling observed with optogenetic stimulation in mouse V1. A.

Schematic of experimental stimulus protocol. If scaling is linear, the same input pulse produces the same response when baseline (spontaneous) rate is changed. **B.** We raise baseline rates using ChR2 in excitatory (E) neurons (Cre-dependent virus in Emx1-Cre mouse line.) Visual stimulus is held constant. **C.** Population histograms showing responses to combined ChR2 and visual (90% contrast) stimuli. Top row: columns show three groups of neurons divided based on size of ChR2 baseline firing rate changes, left: smallest ChR2 effects (N=94; 36 single, 58 multi-units), middle: intermediate ChR2 effects (N=101; 31 single-, 70 multi-units), right: largest ChR2 effects (N=94; 28 single-, 66 multi-units), Brown: responses to visual stimulus with no optogenetic stimulus. Cyan: responses to visual stimulus when baseline rates are changed by sustained optogenetic stimulus. Bottom row: Same data as top row, with spontaneous firing rates subtracted. Visual responses differ

somewhat between columns because each column is a different group of neurons, but within each group there is little response change as spontaneous rate varies. **D**, Linear scaling is seen across a wide contrast range. Top row: responses without baseline subtraction. Bottom row: baseline subtracted. Errorbars: SEM of pooled unit responses. **E**, Linear scaling is seen on average, across neurons with a variety of ChR2-induced baseline rate changes, with some weak sublinearity at the highest rate changes and highest contrasts. Y axes: difference in visual responses (relative to baseline) with and without ChR2 stimulation; dashed line at zero shows a perfectly linear response. Red: lowess regression, shaded region is a bootstrapped 95% confidence interval. Two outlier points in 90% contrast plot are omitted for visual clarity although they are included in the regression; the two outliers are shown in Supp Fig. 1C.

Figure 2: Spiking model shows sublinear scaling with feedforward inhibition.

A, Cartoon of network architecture. Blue: E cells, green: I cells. The conductance-based spiking model produces stochastic V_m and spikes as seen *in vivo*, and an example membrane potential (V_m) trace from one excitatory cell is shown. **B**, Response scaling as feedforward (FF) input to E and I cells is varied. To measure response scaling, inputs to E and/or I cells with specified rate (given by X,Y axes) are delivered, and average response over all E cells is measured. Then, the E and I input rates are multiplied by a constant (here, 2) and the size of the second response is compared to the first. Percent change shown by color, yellow: second response is similar (linear), blue: second response is smaller (sublinear). Contour lines show first response (spk/s). Response rates below 5 spk/s and above 20 spk/s are masked (gray). Average spontaneous rate is adjusted to 5 spk/s (Methods), and 33% of network neurons receive external input, to approximate the sparse set of cortical neurons that typically respond to sensory inputs (Fig. 1). Pink points show E and I rate combinations used in C,D. **C**, Near-linear responses to a range of input sizes when feedforward input is provided to E cells only. Parameters here are indicated by pink dot in B, and first two responses here are the same two responses used to compute percent change shown in color there. Left panel: average rates, right panel: same data replotted showing change (spk/s) in response (y-axis) as a function of prior response (x-axis). For these plots,



a linear response is a horizontal line (dashed gray line). Heavy lines: prior rates less than 50 spk/s, highlighting for visual clarity rates far from saturation caused by absolute refractory period (3 ms). **D**, Sublinear responses to a range of input sizes when input provided to both E and I cells. Same conventions as C. In this case, heavy green line in right panel lies farther below horizontal than heavy blue line in C, showing more sublinear scaling.

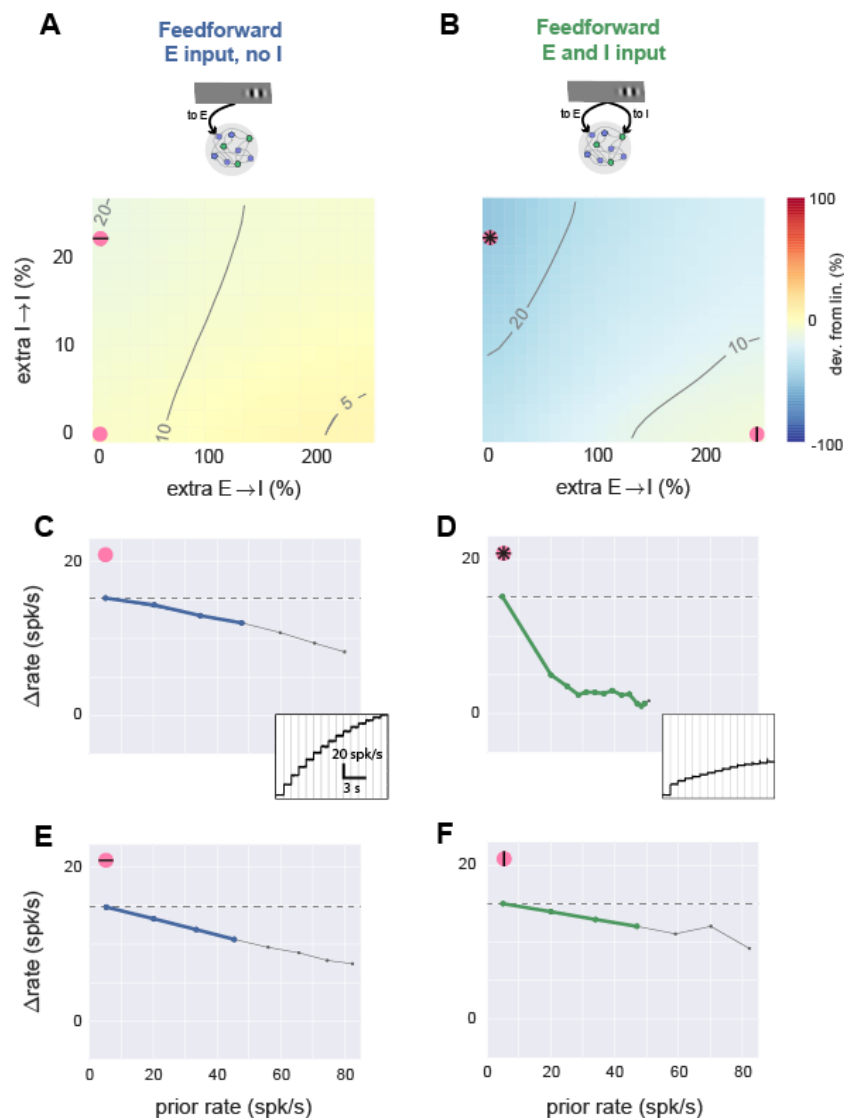


Figure 3: With feedforward inhibition, network model can produce linear or sublinear responses. **A**, Simulations with feedforward input to E cells only, while local network connectivity is varied. X-axis: E to I connection strength, y-axis, I to I connection strength. Axes give percent change in total synaptic input that a single cell receives from one (E or I) population (see Methods), where zero is a balanced network (e.g. Fig. 2) with equal probability of synapses onto E and I cells. Other conventions as in Fig. 2B (contour lines show evoked response to first stimulus, color shows percent difference in response to doubled external stimulus). Spontaneous rate and external stimulus rates are constant for entire panel. **B**, Simulations with feedforward input to E and I cells while local connectivity is varied. Pink symbols show parameter regions where scaling is sublinear (stronger I→I connectivity) or linear (stronger E→I connectivity). **C**, Scaling plot (response size as a function of previous rate) for parameters shown by pink dot in A: no extra local connections, feedforward E only, same parameters as Fig. 2C. Inset: timecourse of responses to the step stimulus; subtracting each rate from rate at the previous step gives y-axis in main panel. **D-F**, same plots, using parameters shown by corresponding pink dots in A-B. Comparing D and E shows that large sublinearity can be produced by extra I→I connections only with feedforward inhibition. Comparing D and F shows that linearity can also be achieved with feedforward inhibition if E→I connectivity is strengthened.

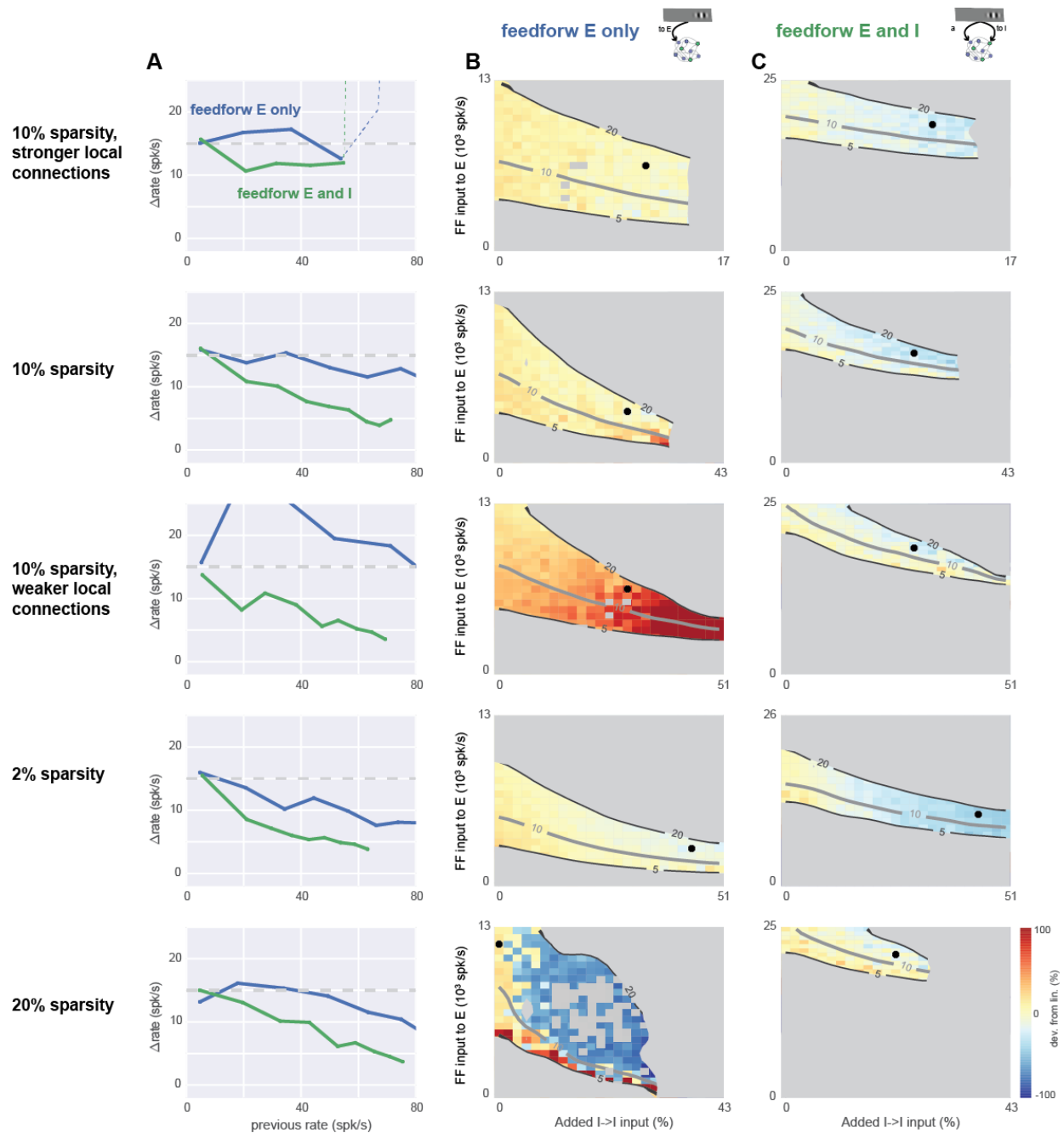


Figure 4: Feedforward inhibition leads to sublinearity in networks with a range of recurrent synaptic sparsities and synaptic strengths. Top row: Simulations in the conductance-based network with 10% connectivity, with strong synapses (each cell receives 10x more E and I input than in the networks of Fig. 2-3). Other rows show networks with different sparsity and synaptic strength. The network of Figs. 2-3 is the fourth row (2% sparsity, 1x strength). **A**, Scaling plots showing network response as a function of prior rate before stimulus. Blue: feedforward E input only but extra I-I input to maximum potential sublinearity, using parameters shown in column B. Green: feedforward E and I input, with extra I-I connections; corresponding parameters are shown in column C. In all rows, feedforward inhibition (green) allows more sublinearity than feedforward excitation alone (blue). Dashed line, top row: network instability (rates diverge). **B**, Average network

response as I->I synaptic strength (x-axis) and feedforward E input (y-axis) are varied. No feedforward inhibition. Black dot shows parameters used to plot blue line in A. Gray regions mask areas where evoked rates are less than 5 spk/s or greater than 20 spk/s, or where network was unstable (rates diverged to maximum rate given by refractory period). Other conventions as in Fig. 2B, 3AB. Feedforward inhibition rate is zero for all rows. **C**, network response as a function of I->I and feedforward E input, in the presence of feedforward inhibition. Individual gray squares seen in fifth row (20% sparsity) column B, inside the 5-20 spk/s contours indicate strongly irregular (non-monotonic) response scaling: strong sublinearity for at least one stimulus step, when both previous and later responses were linear or supralinear. Feedforward inhibition arrival rate to stimulated cells for each row, from top: 14k, 14k, 19k, 11k, 17k spk/s, chosen to give a 15 spk/s response for 3x the feedforward excitatory rate that alone produces a 15 spk/s response (see Fig. 2B). Fourth row (2% sparsity, same network as Fig. 2-3) uses 40% extra I->E connections to show linear responses are robust to many forms of connectivity variation (see also Supp Fig. 1).

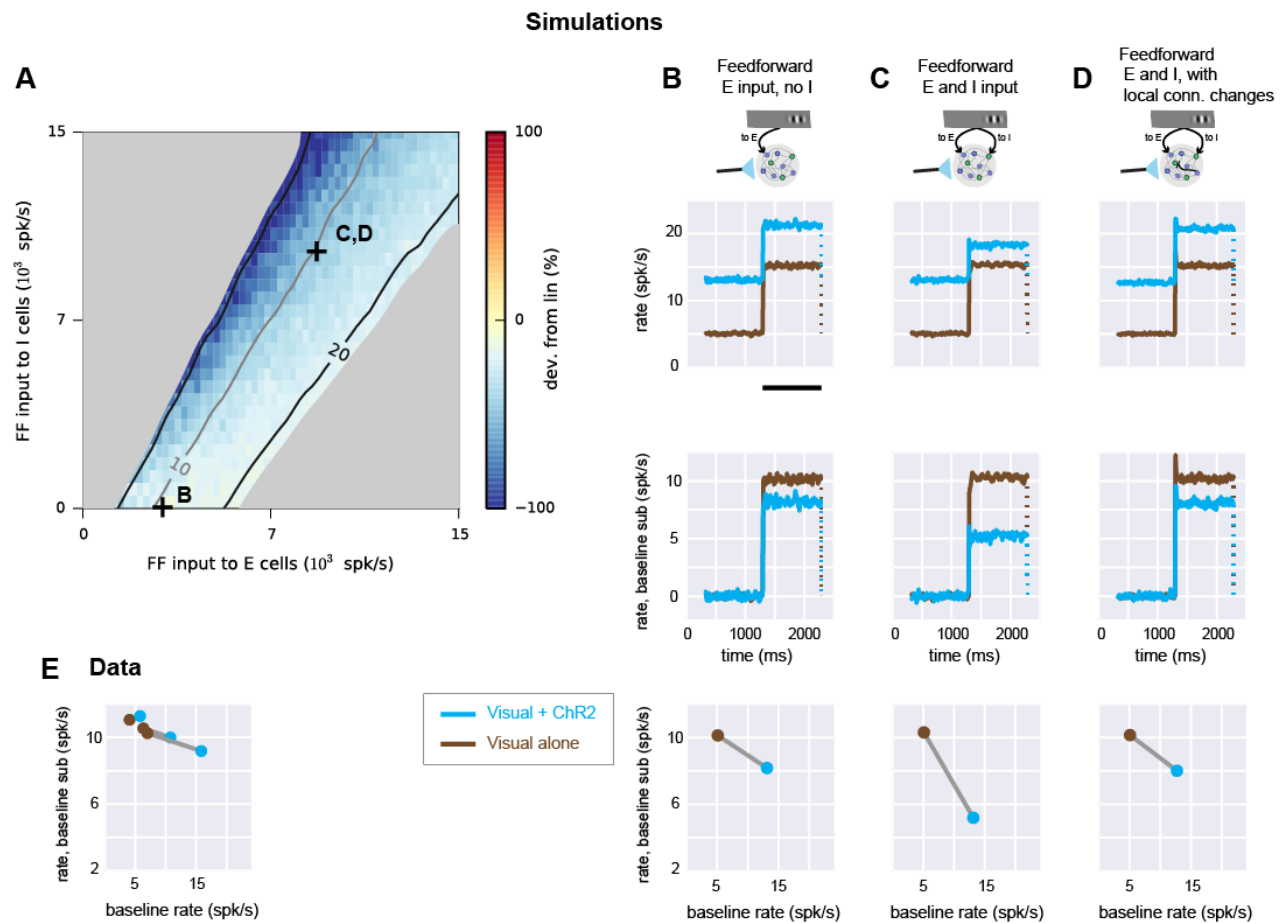
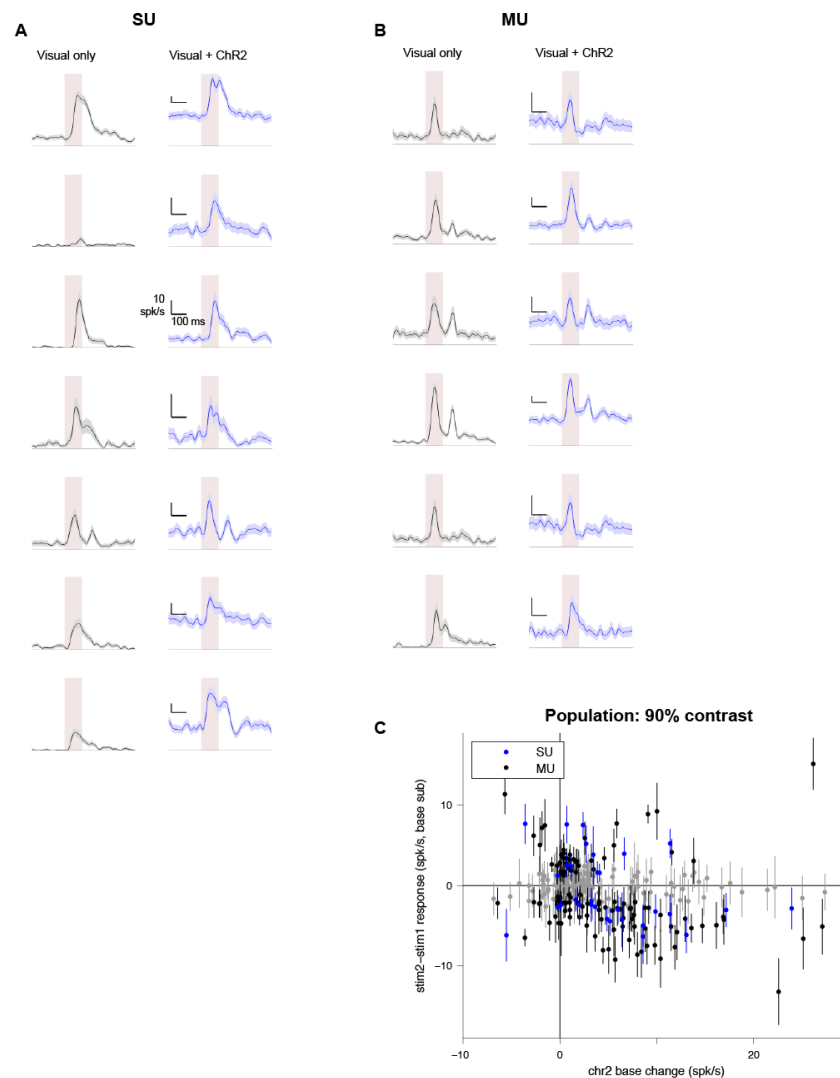


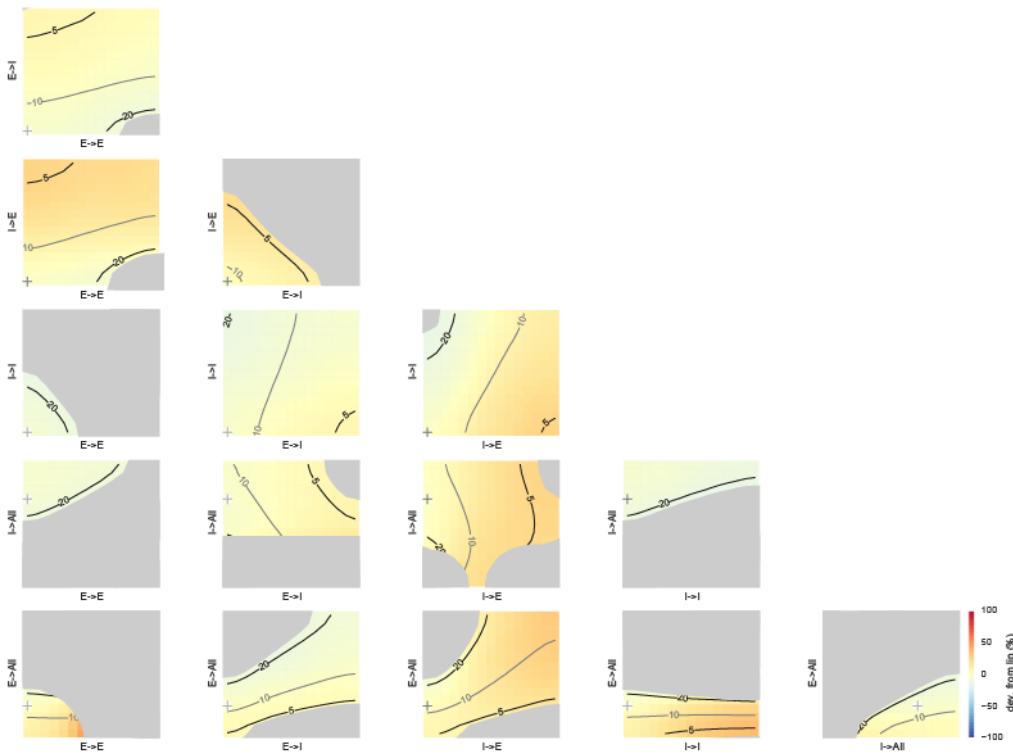
Figure 5: Experimental linear scaling can be replicated in networks receiving feedforward inhibition. **A**, Simulation where conductance steps (modeling ChR2 input) and feedforward Poisson trains (modeling visual input) are combined. Strength of feedforward E input (x-axis) and feedforward I input (y-axis) are varied while spontaneous rate is set to 5 spk/s. Connection sparsity is 2%. Other conventions as in Fig. 2B. (+) symbols show values of E,I input used in panels B-D. **B**, network responses when feedforward input is supplied to E cells only. Top row: network responses (mean of E cell rates). Brown: feedforward Poisson (visual) input only. Cyan: conductance (ChR2) input combined with visual input. Conductance increase lasts for the full duration of the cyan trace. Visual input duration is shown by black bar (bottom of plot). Dotted line indicates rates return to previous baseline when feedforward input ends. Second row: same data as top row, with baseline rate subtracted. Third row: response (y-axis) as a function of rate before feedforward input begins (x-axis). **C**, same network simulations with feedforward input to both E and I cells (parameters marked by C in panel A). **D**, network receiving feedforward input to both E and I cells, but with stronger local connections from E to I cells (cf. Fig. 3, similar effect for two feedforward Poisson inputs, instead of feedforward input paired with conductance step as shown here). **E**, data from Fig. 1C plotted to show how responses scale as baseline is changed. Three lines (brown: no ChR2, cyan: with ChR2) are the three groups of recorded neurons shown in Fig. 1C.

Supplementary Figures

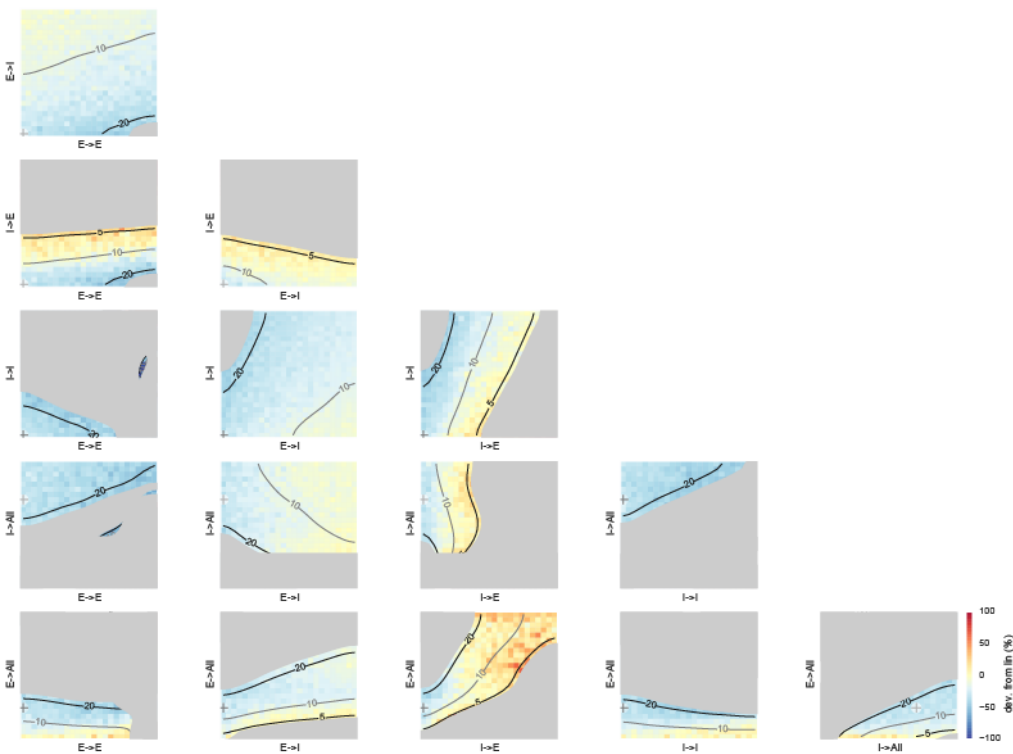


Supplementary Figure 1: Individual unit responses show both supra- and sub-linear responses, though mean response is nearly linear. **A**, Pairs of plots in each row show data from individual well-isolated example units (single units), all from Emx1-ChR2 stimulation of excitatory neurons. Left panels: average response to a visual stimulus without ChR2 stimulation (shaded areas, \pm SEM, across ~ 100 stimulus repetitions). Right panels: average response to visual stimulus with ChR2 stimulation. **B**, same as **A**, showing example multi-units (waveforms not well isolated from other cells; Methods). **C**, Population of recorded units, showing that many individual units are significantly supra- or sub-linear. Visual stimulus has 90% contrast. X-axis: average firing rate change with ChR2 stimulus, Y-axis: change in visual responses (each visual response measured from baseline) with and without optogenetic stimulus. Errorbars: SEM. Points that are at least 1 SEM away from horizontal line at zero (linear response) are colored blue (SU) or black (MU). Points within 1 SEM of linear are colored gray. Data are as in Fig. 1E for 90% contrast, here with std. err. for each point and adding on the negative Y-axis the few units that are suppressed by stimulation. 34% of single units are significantly non-linear (17/50, $p < 0.01$, KS test), and 28% of multi-units are significantly non-linear (67/239, $p < 0.01$, KS test).

A Feedforward input to E only - close to linear



B Feedforward input to E and I cells - linear or sublinear



Supplementary Figure 2: Near-linear scaling is seen even when connectivity is varied, unless feedforward inhibition is present. **A**, no feedforward inhibition. The network is the

2% sparsity network shown in Fig. 2B-D, Fig. 3, Fig. 4D. Each panel is a pairwise combination of synaptic strength parameters. The six strength parameters are: extra connections between E and I populations (E->E, E->I, I->I, I->I), plus the standard sparse connections from both populations to all other cells in the network (labeled E->All, I->All). As in pairwise color plots in Figs. 3AB, 4BC, x-axis and y-axis represent strength of all the synapses in that set of connections (one of the six groups of synapses listed above). Conventions as in Fig. 3AB: Axes show variation in synapse strength, contour lines show evoked response to fixed-size external input step, colors show deviation from linear scaling when response to an external input step is compared to response to an input step scaled by 2x. Spontaneous rate throughout is set at 5 spk/s (Methods). Feedforward external input strength is same as marked by open circle (O) in Fig. 2B. (+) symbol in all panels of this figure shows parameters used for Fig. 2B-D where all extra parameters are zero. **B**, Both feedforward excitation and feedforward inhibition. Same conventions as A. Feedforward external input strength is same as marked by (X) symbol in Fig. 2. (+) symbol in this figure shows parameters used for Fig. 2B-D (*i.e.* with all extra connection strength parameters zero). Several regions in the parameter space show sublinear scaling (blue).

References

- Ahmadian, Y., Rubin, D.B., and Miller, K.D. (2013). Analysis of the stabilized supralinear network. *Neural computation* 25, 1994-2037.
- Barral, J., and Reyes, A.D. (2016). Synaptic scaling rule preserves excitatory-inhibitory balance and salient neuronal network dynamics. *Nature Neuroscience* 19, 1690-1700.
- Bonin, V., Histed, M.H., Yurgenson, S., and Reid, R.C. (2011). Local diversity and fine-scale organization of receptive fields in mouse visual cortex. *Journal of Neuroscience* 31, 18506-18521.
- Bonin, V., Mante, V., and Carandini, M. (2006). The statistical computation underlying contrast gain control. *Journal of Neuroscience* 26, 6346-6353.
- Braitenberg, V., and Schüz, A. (2001). Cortex: Statistics and Geometry of Neuronal Connectivity. In: *Neuro-und Sinnesphysiologie*.
- Brette, R., Rudolph, M., Carnevale, T., Hines, M., Beeman, D., Bower, J.M., Diesmann, M., Morrison, A., Goodman, P.H., Harris, F.C., *et al.* (2007). Simulation of networks of spiking neurons: a review of tools and strategies. *Journal of computational neuroscience* 23, 349-398.
- Briggman, K.L., Helmstaedter, M., and Denk, W. (2011). Wiring specificity in the direction-selectivity circuit of the retina. *Nature* 471, 183-188.
- Brunel, N. (2000). Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons. *Journal of computational neuroscience* 8, 183-208.
- Carandini, M., and Heeger, D.J. (2012). Normalization as a canonical neural computation. *Nat Rev Neurosci* 13, 51-62.
- Chance, F.S., Abbott, L.F., and Reyes, A.D. (2002). Gain modulation from background synaptic input. *Neuron* 35, 773-782.
- Chen, J.L., Flanders, G.H., Lee, W.-C.A., Lin, W.C., and Nedivi, E. (2011). Inhibitory dendrite dynamics as a general feature of the adult cortical microcircuit. *Journal of Neuroscience* 31, 12437-12443.
- Connors, B.W., Gutnick, M.J., and Prince, D.A. (1982). Electrophysiological properties of neocortical neurons in vitro. *J Neurophysiol* 48, 1302-1320.
- Cruikshank, S.J., Lewis, T.J., and Connors, B.W. (2007). Synaptic basis for intense thalamocortical activation of feedforward inhibitory cells in neocortex. *Nat Neurosci* 10, 462-468.
- Destexhe, A., and Paré, D. (1999). Impact of network activity on the integrative properties of neocortical pyramidal neurons in vivo. *J Neurophysiol* 81, 1531-1547.

- Destexhe, A., Rudolph, M., and Paré, D. (2003). The high-conductance state of neocortical neurons in vivo. *Nat Rev Neurosci* 4, 739-751.
- Douglas, R.J., Koch, C., Mahowald, M., Martin, K.A., and Suarez, H.H. (1995). Recurrent excitation in neocortical circuits. *Science* 269, 981-985.
- Douglas, R.J., and Martin, K.A.C. (2004). Neuronal circuits of the neocortex. *Annu Rev Neurosci* 27, 419-451.
- Freund, T.F., Martin, K.A., Soltesz, I., Somogyi, P., and Whitteridge, D. (1989). Arborisation pattern and postsynaptic targets of physiologically identified thalamocortical afferents in striate cortex of the macaque monkey. *J Comp Neurol* 289, 315-336.
- Gabernet, L., Jadhav, S.P., Feldman, D.E., Carandini, M., and Scanziani, M. (2005). Somatosensory integration controlled by dynamic thalamocortical feed-forward inhibition. *Neuron* 48, 315-327.
- Garraghty, P.E., and Sur, M. (1990). Morphology of single intracellularly stained axons terminating in area 3b of macaque monkeys. *J Comp Neurol* 294, 583-593.
- Glickfeld, L.L., Histed, M.H., and Maunsell, J.H.R. (2013). Mouse primary visual cortex is used to detect both orientation and contrast changes. *J Neurosci* 33, 19416-19422.
- Gorski, J.A., Talley, T., Qiu, M., Puelles, L., Rubenstein, J.L.R., and Jones, K.R. (2002). Cortical excitatory neurons and glia, but not GABAergic neurons, are produced in the Emx1-expressing lineage. *J Neurosci* 22, 6309-6314.
- Häusser, M., Spruston, N., and Stuart, G.J. (2000). Diversity and dynamics of dendritic signaling. *Science* 290, 739-744.
- Histed, M.H., and Maunsell, J.H.R. (2013). Cortical neural populations can guide behavior by integrating inputs linearly, independent of synchrony. *Proc Natl Acad Sci USA*.
- Huang, X., Elyada, Y.M., Bosking, W.H., Walker, T., and Fitzpatrick, D. (2014). Optogenetic assessment of horizontal interactions in primary visual cortex. *Journal of Neuroscience* 34, 4976-4990.
- Isaacson, J.S., and Scanziani, M. (2011). How inhibition shapes cortical activity. *Neuron* 72, 231-243.
- Katz, Y., Menon, V., Nicholson, D.A., Geinisman, Y., Kath, W.L., and Spruston, N. (2009). Synapse distribution suggests a two-stage model of dendritic integration in CA1 pyramidal neurons. *Neuron* 63, 171-177.
- Koike, H., Mano, N., Okada, Y., and Oshima, T. (1970). Repetitive impulses generated in fast and slow pyramidal tract cells by intracellularly applied current steps. *Exp Brain Res* 11, 263-281.

- Lee, W.-C.A., Bonin, V., Reed, M., Graham, B.J., Hood, G., Glattfelder, K., and Reid, R.C. (2016). Anatomy and function of an excitatory network in the visual cortex. *Nature* 532, 370-374.
- Lefort, S., Tómm, C., Floyd Sarria, J.-C., and Petersen, C.C.H. (2009). The excitatory neuronal network of the C2 barrel column in mouse primary somatosensory cortex. *Neuron* 61, 301-316.
- Li, Y.-t., Ibrahim, L.A., Liu, B.-h., Zhang, L.I., and Tao, H.W. (2013). Linear transformation of thalamocortical input by intracortical excitation. *Nat Neurosci* 16, 1324-1330.
- Lien, A.D., and Scanziani, M. (2013). Tuned thalamic excitation is amplified by visual cortical circuits. *Nat Neurosci* 16, 1315-1323.
- Nassi, J.J., Avery, M.C., Cetin, A.H., Roe, A.W., and Reynolds, J.H. (2015). Optogenetic Activation of Normalization in Alert Macaque Visual Cortex. *Neuron* 86, 1504-1517.
- Nikolic, K., Grossman, N., Grubb, M.S., Burrone, J., Toumazou, C., and Degenaar, P. (2009). Photocycles of channelrhodopsin-2. *Photochem Photobiol* 85, 400-411.
- Olsen, S.R., Bortone, D.S., Adesnik, H., and Scanziani, M. (2012). Gain control by layer six in cortical circuits of vision. *Nature* 483, 47-52.
- Priebe, N.J., and Ferster, D. (2008). Inhibition, spike threshold, and stimulus selectivity in primary visual cortex. *Neuron* 57, 482-497.
- Reid, R.C. (2001). Divergence and reconvergence: multielectrode analysis of feedforward connections in the visual system. *Prog Brain Res* 130, 141-154.
- Renart, A., de la Rocha, J., Bartho, P., Hollender, L., Parga, N., Reyes, A., and Harris, K.D. (2010). The asynchronous state in cortical circuits. *Science* 327, 587-590.
- Richardson, M.J.E. (2004). Effects of synaptic conductance on the voltage distribution and firing rate of spiking neurons. *Phys Rev E Stat Nonlin Soft Matter Phys* 69, 051918.
- Richardson, M.J.E. (2007). Firing-rate response of linear and nonlinear integrate-and-fire neurons to modulated current-based and conductance-based synaptic drive. *Phys Rev E Stat Nonlin Soft Matter Phys* 76, 021919.
- Rubin, D.B., Van Hooser, S.D., and Miller, K.D. (2015). The stabilized supralinear network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron* 85, 402-417.
- Sato, T.K., Häusser, M., and Carandini, M. (2014). Distal connectivity causes summation and division across mouse visual cortex. *Nat Neurosci* 17, 30-32.
- Silver, R.A. (2010). Neuronal arithmetic. *Nat Rev Neurosci* 11, 474-489.

- Steriade, M., Timofeev, I., and Grenier, F. (2001). Natural waking and sleep states: a view from inside neocortical neurons. *J Neurophysiol* 85, 1969-1985.
- Swadlow, H.A. (2003). Fast-spike interneurons and feedforward inhibition in awake sensory neocortex. *Cereb Cortex* 13, 25-32.
- Tiesinga, P., Fellous, J.-M., and Sejnowski, T.J. (2008). Regulation of spike timing in visual cortical circuits. *Nat Rev Neurosci* 9, 97-107.
- Van Essen, D.C., Newsome, W.T., and Maunsell, J.H. (1984). The visual field representation in striate cortex of the macaque monkey: asymmetries, anisotropies, and individual variability. *Vision Res* 24, 429-448.
- van Vreeswijk, C., and Sompolinsky, H. (1996). Chaos in neuronal networks with balanced excitatory and inhibitory activity. *Science* 274, 1724-1726.
- van Vreeswijk, C., and Sompolinsky, H. (1998). Chaotic balanced state in a model of cortical circuits. *Neural computation* 10, 1321-1371.
- Vogels, T.P., and Abbott, L.F. (2005). Signal propagation and logic gating in networks of integrate-and-fire neurons. *Journal of Neuroscience* 25, 10786-10795.
- Yamins, D.L.K., and DiCarlo, J.J. (2016). Using goal-driven deep learning models to understand sensory cortex. *Nat Neurosci* 19, 356-365.