# Learning quadratic receptive fields from neural responses to natural stimuli

Kanaka Rajan[*a], Olivier Marre[b] and Gašper Tkačik[†c]

[a] Joseph Henry Laboratories of Physics,
Lewis–Sigler Institute for Integrative Genomics, Princeton University, Princeton NJ 08544, USA
[b] Institution de la Vision, UPMC UMRS 968, INSERM, CNRS U7210, CHNO Quinze-Vingts, F-75012 Paris, France
[c] Institute of Science and Technology Austria, Am Campus 1, A-3400 Klosterneuburg, Austria

Models of neural responses to stimuli with complex spatiotemporal correlation structure often assume that neurons are only selective for a small number of linear projections of a potentially high-dimensional input. Here we explore recent modeling approaches where the neural response depends on the quadratic form of the input rather than on its linear projection, that is, the neuron is sensitive to the local covariance structure of the signal preceding the spike. To infer this quadratic dependence in the presence of arbitrary (e.g. naturalistic) stimulus distribution, we review several inference methods, focussing in particular on two information-theory-based approaches (maximization of stimulus energy or of noise entropy) and a likelihood-based approach (Bayesian spike-triggered covariance, extensions of generalized linear models). We analyze the formal connection between the likelihood-based and information-based approaches to show how they lead to consistent inference. We demonstrate the practical feasibility of these procedures by using model neurons responding to a flickering variance stimulus.

## I. INTRODUCTION

A basic challenge in sensory neuroscience has been to develop a mathematically concise description of how neurons encode stimuli into sequences of spikes. There are two main approaches to this task, which differ primarily in how much emphasis is placed on anatomical structure versus function. Structure-based modeling starts at the level where basic physical processes govern the observed phenomena. A realistic, conductance-based model could thus be used to predict the neuron's response to a particular type of applied stimulation (Hodgkin & Huxley, 1952; Koch, 1999). While this bottom-up approach is directly interpretable in terms of biophysical components and processes, it has a number of disadvantages: (i) the required parameters might be experimentally inaccessible; (ii) in a sensory context, the inputs in this model (the activity of presynaptic neurons) could be related in a complex and intractable manner to the stimulus under experimental control; and, (iii), with enough modeling detail, the problem of understanding or summarizing the "computation" that the model implements can become as difficult as understanding the real neuron itself.

Functional models, in contrast to the above, attempt to capture only the essence of the neural computation: the transformation of stimuli into spiking responses (see Wu et al. (2006) for an in-depth review). These models are usually fully learned from data, rather than being derived from the underlying dynamical or physical model (but see Agüera y Arcas et al. (2003); Agüera y Arcas & Fairhall (2003); Hong et al. (2007); Lundstrom et al. (2008); Ostrojic & Brunel (2011)). Two considerations are therefore critical to the success of functional models: whether typical electrophysiological recordings can provide enough data for successful inference of the model's parameters; and whether efficient inference algorithms for these parameters exist. Because the space of all possible stimuli (e.g. all images incident on a retina) and all possible responses (e.g. complete sets of spike arrival times) is vast, our progress must depend on making well-chosen simplifying assumptions. One extreme simplification, for example, involves varying the stimulus along one "dimension" only, as in the case of the orientation or wavelength of a drifting grating visual stimulus, and representing the output by a single scalar quantity, e.g. the average firing rate in a chosen time bin. These measurements have traditionally been summarized by tuning curves and have provided basic insights into principles of sensory and population coding (Dayan & Abbott, 2001). The relevance of the tuning curve approach is, however, limited by the choice of the single dimension along which the stimulus is manipulated, which may drastically underestimate the true complexity in the structure of the stimuli to which the neuron could respond. Despite strong limitations, such studies helped establish the concept of a "receptive field," the region of stimulus space where changes in the stimulus modulate the spiking behavior of the neuron.

Central to the concept of receptive field is the notion of locality in the stimulus or feature space. For instance, a ganglion cell in the retina may be sensitive only to specific changes in light intensity that occur within a small visual
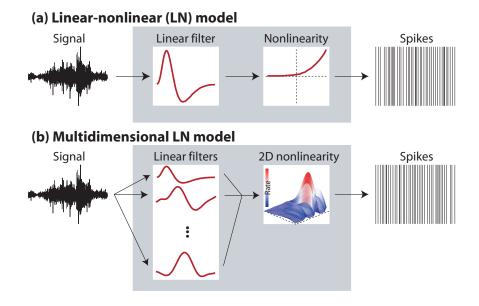
[*] krajan@princeton.edu

[†] gtkacik@ist.ac.at

arXiv:1209.0121v1 [q-bio.NC] 1 Sep 2012

**(a) Linear-nonlinear (LN) model**



**(b) Multidimensional LN model**



FIG. 1 **A schematic showing linear-nonlinear architectures.** **(a)** The instantaneous firing rate, or probability per unit time of emitting a spike, in a linear-nonlinear model neuron is obtained by passing the signal through a linear filter, and mapping the resulting value through a pointwise nonlinearity. **(b)** A multidimensional LN model neuron requires the signal to be filtered through $K$ linear filters. The number of filters, $K$, is usually much smaller than the dimension of the stimulus. The stimulus projections are mapped into the firing rate through a $K$-dimensional nonlinear function. Without further assumptions, inference of these models is tractable from real recordings only when $K$ is small (usually less than 3).

angle (Hartline, 1940). A productive way of capturing this notion of locality has been to think of a receptive field as one or more filters that act on the stimulus; only those stimulus variations that result in the change in filter output have the ability to affect the neural response. In this view, the neurons perform dimensionality reduction by projecting the stimulus down into a small number of dimensions. Consequently, the success of data analysis techniques built around this idea must depend on whether a small number of filters suffices to fully account for the neuron's sensitivity and its response properties.

Methods based in systems identification theory have provided systematic procedures to infer both the receptive fields of neurons as well as subsequent computations. These methods usually share two key features. First, they can (sometimes necessarily) be used with stimuli that sample the stimulus space broadly, making no explicit assumptions about which stimulus features are important. This is in contrast to the restricted stimuli employed for measuring tuning curves. Second, the procedures usually involve a series of approximations that can provably yield an ever better description of the system if increasing amounts of data are available. Table I provides an overview of various functional models and related inference methods. Among the earliest to be used successfully, Wiener and Volterra expansions helped identify the first- and second-order kernels mapping the stimulus to response time traces in various systems (Marmarelis & Marmarelis, 1978; Recio-Spinoso et al., 2005; Sakai, 1992; Schetzen, 1989; Victor & Knight, 1979; Wiener, 1958). However, in many cases strong intrinsic nonlinearities attributed to spike generation would require a large number of terms in Wiener-Volterra expansions, despite the fact that the underlying stimulus sensitivity might be simpler, and therefore of low order. Models where the (possibly linear) projections of the stimulus in the receptive field were decoupled from the nonlinearities underlying spiking, as in linear-nonlinear (LN) architectures illustrated in Fig. 1, made further progress possible.

LN and LN-like models have been used widely and profitably to predict the firing rate traces of single sensory neurons, because their parameters can be inferred easily under suitable conditions. However, the more intriguing cases are the ones where LN models either perform poorly or fail entirely. One such failure mode is the inability to account for the statistics of neural activity beyond the mean firing rate. Specifically, real sensory neurons often have variability that is smaller than that attributed to Poisson processes (de Ruyter van Steveninck et al., 1997); phenomena like refractoriness and spike rate adaptation are not captured by LN models (Berry & Meister, 1998); and in neural populations, uncoupled LN models fail to reproduce the basic covariance structure of neural activity (Granot-Atedgi et al., 2012; Pillow et al., 2008; Schneidman et al., 2006). Some of these issues can be addressed by adding suitable dynamical complexity beyond the linear filtering stage, to make the nonlinearities in spike generation more realistic (Keat et al., 2001; Ozuysal & Baccus, 2012; Paninski et al., 2004), or by including interactions between neurons in models of neural firing (Granot-Atedgi et al., 2012; Pillow et al., 2008).

A different kind of failure of LN models rests on the assumption that stimulus sensitivity occurs through a single (or a small number of) linear projections. One example is contrast adaptation, where a simple LN model derived from a white noise stimulus of a certain variance fails to accurately predict the response to a stimulus with smaller or larger variance (Baccus & Meister, 2002; Borst & Egelhaaf, 1987; van Hateren, 1992; de Ruyter van Steveninck et al., 1986; Smirnakis et al., 1997). Other examples include the failure to account for the sensitivity of retinal ganglion cells to fine spatial detail (possibly because of nonlinear summation within the receptive field (Demb et al., 2001)), or to stimulus motion (Berry et al., 1999; Gollisch & Meister, 2010; Schwartz et al., 2007). Generally, these difficulties emerge clearly when the stimulus statistics change or increase in complexity beyond those used to infer the model, for instance by becoming more "naturalistic"– i.e. having temporal and spatial pairwise correlation over many scales, skewed first order histograms, and statistical structure beyond second order.

The problems with the LN models can generally be addressed in two possible ways. In the first, LN models can be extended to account for a particular phenomenon on a particular stimulus, e.g., by adding a contrast gain control mechanism (Schwartz & Simoncelli, 2001; Schwartz et al., 2002) or by an *ad hoc* rescaling of nonlinearities (Brenner et al., 2000) to account for contrast adaptation in an experiment where the variance of a Gaussian input is modulated. The second approach is more general: by using complex stimuli, including fully natural movies from the start, the goal is to find the complete (or close to complete) set of features to which the neuron responds. It is worth noting that these two approaches, as well as the associated use of simple vs natural stimulus ensembles, generally reflect two motivations for building models of neural encoding in the first place: one is to propose and test specific simple models and incrementally improve them, while the other is to infer descriptions that should be valid across a wide range of stimuli and conditions from the start. For the former purpose—falsifying a model or developing a simple functional form for the stimulus-response relationship—using a stimulus set that is analytically convenient but highly un-natural, e.g., white noise, is sufficient. This is because when a proposed model fails on a subset of stimuli, it can be excluded or must be extended by additional mechanisms. Until recently, this was the main reason for using systems identification methods with noise stimuli. The drive to use naturalistic stimuli comes, on the other hand, from trying to find a model that captures from the start the responses to a wide variety of biologically relevant inputs, and from the observation that naturalistic stimuli may change even the basic filter responses of cells (Sharpee et al., 2006) and engage response mechanisms that are difficult to probe using noise stimulation (e.g., Olveczky et al. (2007)). Potential drawbacks with using natural stimuli include technical obstacles in model inference and the statistical intractability of the natural ensemble (Geisler, 2008; Simoncelli & Olshausen, 2001). The choice of the stimulus ensemble certainly deserves a lengthier discussion; see, for example, Rust & Movshon (2005).

To find the complete set of stimulus features to which a neuron responds, one can look for multiple linear features, a task for which methodological frameworks exist and have been validated for a small number of features. Unfortunately, extracting more than 2 or 3 features becomes intractable because of the curse of dimensionality. A possible anatomically motivated simplification of a multi-feature LN model is a cascade LN (an LNLN) model, where the nonlinearly transformed filter outputs are linearly summed and passed through the spike-generating nonlinearity. Despite some successes (Bölinger & Gollisch, 2012; Gollisch & Herz, 2005), the general problem of inferring cascading models remains technically challenging (usually involving difficult optimizations). A somewhat simpler LNL system has proven both to account for the behavior of the Y-type retinal ganglion cells very well, as well as being tractable to infer using the sum-of-sinusoids formulation of the Wiener formalism (Victor & Knight, 1979; Victor & Shapley, 1979, 1980). A particular case of interest for this review is a special subclass of LNLN models which can be reformulated as quadratic-nonlinear models, i.e. models where the initial dimensionality reduction of the stimulus is not a linear projection of the stimulus, but rather an arbitrary quadratic function of the stimulus.

Recently there has been a lot of interest in designing systematic, tractable methods for inferring neural sensitivities when the initial dimensionality reduction step is of high-order (e.g. quadratic)[1]. In this paper, we start by presenting several biologically motivated examples of quadratic stimulus sensitivity in Section II. We then review several complementary approaches that can be used to learn quadratic stimulus dependence even when neurons are responding to rich, naturalistic stimuli: we discuss the maximally informative stimulus energy (Rajan & Bialek, 2012) and the maximization of noise entropy (Fitzgerald et al., 2011a,b; Globerson et al., 2009) in Section III.A, and follow with the Bayesian spike-triggered covariance (Park & Pillow, 2011) and related extensions of generalized linear models to quadratic stimulus dependence[2] in Section III.C. We show under which conditions information and likelihood based approaches lead to consistent inference in the Appendix.

––––––

[1] When we speak of the order (e.g. linear, quadratic etc), we refer to the order of the kernel operating on the stimulus, which can be defined unambiguously. In contrast, the order of the neural processing system as a whole depends on the stimulus statistics; for example, higher-order statistical structure in the stimulus can conflate first- and second-order responses of the system. Likewise, aspects of the response explained by second order kernel inferred even with Gaussian noise depend on the power spectrum of the input.

[2] This problem has been worked on by the authors of this review in parallel with the authors of Park & Pillow (2011).

| Method | Stimulus type | Models / restrictions | References |
|---|---|---|---|
| Wiener/Volterra series | white gaussian noise, sum-of-sinusoids | $r = r_0 + \mathbf{k} \cdot \mathbf{s} + \mathbf{s}^T \mathbf{Q}\mathbf{s} + \cdots$ | (Marmarelis & Marmarelis, 1978; Schetzen, 1989; Wiener, 1958) (Recio-Spinoso et al., 2005; Victor & Knight, 1979) |
| spike trigger average (STA) (reverse correlation) | spherically symmetric, binary noise, m-sequences | LN (single filter), isolated spikes $r = f(\mathbf{k} \cdot \mathbf{s})$ | (de Boer & Kuyper, 1968; Paninski, 2003; Reid et al., 1997) (Schwartz et al., 2006; Simoncelli et al., 2004) |
| debiased STA (reverse correlation) | "gaussian-like" asym. 1-point histogram | LN (single filter), isolated spikes $r = f(\mathbf{k} \cdot \mathbf{s})$ | (Lesica et al., 2008) |
| spike trigger covariance (STC) (reverse correlation) | gaussian | LN (multiple filters), isolated spikes $r = f(\mathbf{k}_1 \cdot \mathbf{s}, \ldots, \mathbf{k}_K \cdot \mathbf{s})$ | (Bialek & de Ruyter van Steveninck, 2005; de Ruyter van Steveninck & Bialek, 1988) (Fairhall et al., 2006; Maravall et al., 2007; Schwartz et al., 2002; Simoncelli et al., 2004) |
| extended projection pursuit regression (ePPR) | any | LN (multiple filters) $r = f(\mathbf{k}_1 \cdot \mathbf{s}, \ldots, \mathbf{k}_K \cdot \mathbf{s})$ | (Rapela et al., 2010) |
| iSTAC (reverse correlation) | gaussian | LN (multiple filters) $r = f(\mathbf{k}_1 \cdot \mathbf{s}, \ldots, \mathbf{k}_K \cdot \mathbf{s})$ | (Pillow & Simoncelli, 2006) |
| differential reverse correlation (dRC) (reverse correlation) | spike triggering snippet | linear feature that predicts spike timing $t_{spike} \propto \mathbf{k} \cdot \mathbf{s}$ | (Tkačik & Magnasco, 2008) |
| maximally informative dimensions (MID) (info maximization) | any | LN (multiple filters) $r = f(\mathbf{k}_1 \cdot \mathbf{s}, \ldots, \mathbf{k}_K \cdot \mathbf{s})$ | (Kouh & Sharpee, 2009; Sharpee et al., 2004, 2006) |
| (maximum likelihood) | any | leaky integrate and fire (LIF/LN-LIF) | (Gerstner & Kistler, 2002; Paninski et al., 2004; Pillow, 2007) |
| error function minimization (general fitting methods) | any | dynamical extensions of LN $r = \Theta(h)\dot{h}, h = \mathbf{k} \cdot \mathbf{s} + \mathbf{q} \cdot \mathbf{y} + \eta$ (Keat), $\dot{A}_i = M_{ij}(f(\mathbf{k} \cdot \mathbf{s}))A_j, r = A_1$ (LNK) | (Keat et al., 2001; Ozuysal & Baccus, 2012) |
| generalized linear models (GLM) (maximum likelihood) | any | point process (dependence on past spiking) $r = f(\mathbf{k} \cdot \mathbf{s} + \mathbf{q} \cdot \mathbf{y} +$ (effect of other neurons) ) | (Paninski, 2004; Pillow et al., 2008; Truccolo et al., 2004) (Gerwinn et al., 2010; Pillow, 2007) |
| isoresponse mapping | synthetic stimuli (parametrizable, low-D) | LNLN cascade $r = f(\mathbf{k}_1 * g(\mathbf{k}_2 * \mathbf{s}))$ | (Bölinger & Gollisch, 2012; Gollisch & Herz, 2005) |
| maximally informative stim. energy (MISE) (info maximization) | any | general quadratic model $r = f(\mathbf{k} \cdot \mathbf{s}, \mathbf{s}^T \mathbf{Q}\mathbf{s})$ | (Rajan & Bialek, 2012) |
| maximization of noise entropy (convex optimization) | any | $r = \text{logistic}(k_0 + \mathbf{k} \cdot \mathbf{s} + \mathbf{s}^T \mathbf{Q}\mathbf{s})$ | (Fitzgerald et al., 2011a,b; Globerson et al., 2009) |
| Bayesian STC / quadratic GLM (likelihood maximization) | any | additive linear and quadratic contributions $r = f(\mathbf{k} \cdot \mathbf{s} + \mathbf{s}^T \mathbf{Q}\mathbf{s})$ | (Park & Pillow, 2011) |

TABLE I **Functional models for single neurons and the related inference methods.** $r(t)$ is the firing rate or the probability of spiking; $\mathbf{k}$ are linear filters acting on stimulus clips $\mathbf{s}$; $\mathbf{Q}$ is a quadratic kernel (any symmetric matrix); $\mathbf{q}$ is a linear filter on the sequence of past spikes $\mathbf{y}$; $f, g$ are arbitrary nonlinear functions; $*$ denotes a convolution; $\Theta$ is a thresholding operation (1 when the argument crosses some threshold from below, 0 otherwise); $\eta$ is a white noise Langevin force. In this paper, we use the term "gaussian" to denote stimuli whose components are jointly Gaussian and possibly correlated (i.e. non-white), unless otherwise stated. While reverse correlation methods are formally simpler for uncorrelated (white) Gaussian noise, it is possible to generalize them for use with correlated noise ensembles. For example, to compute an unbiased estimate of the linear (L) part of the model using STA and a correlated stimulus, one needs to correct for stimulus correlations by acting on the spike triggered average with the inverse covariance matrix. For an extensive review of spike-triggered (reverse correlation) methods, see Schwartz et al. (2006).

## II. HIGH-ORDER STIMULUS DEPENDENCE

In a typical experiment, a neuron can be driven by a synthetically generated stimulus containing a desired statistical structure. For probing the visual system for example, this stimulus might be a random binary checkerboard, a drifting grating, or full-field light intensity flicker. If the neuron's response depends solely on the stimulus presented in the recent past of duration $T$ (and possibly on its own previous spiking behavior), we can restrict our attention to stimulus clips $\mathbf{s}$ of length $\geq T$. These clips are drawn from a distribution $P(\mathbf{s})$ that characterizes the stimulus; the $N$ components of vector $\mathbf{s}$ represent successive stimulus values in time and optionally across space. Our task is then to infer the dependence of the instantaneous probability of spiking (firing rate) at time $t$ on the stimulus, $\mathbf{s}(t)$, presented just prior to $t$.

If the neuron is well described by the linear-nonlinear (LN) model, where the spiking rate $r$ is an arbitrary positive, point-wise, nonlinear function $f$ of the stimulus projected onto the filter, $r(\mathbf{s}) = f(\mathbf{k} \cdot \mathbf{s})$, and the stimulus distribution is chosen to be spherically symmetric, $P(\mathbf{s}) = P(|\mathbf{s}|)$, we can use the spike-triggered average (STA) to obtain an unbiased estimate of the single linear filter $\mathbf{k}$ (de Boer & Kuyper, 1968; Simoncelli et al., 2004). Spike-triggered covariance (STC) generalizes the filter inference to cases where the firing rate depends nonlinearly on $K \geq 1$ projections of the stimulus, $r(\mathbf{s}) = f(\mathbf{k}_1 \cdot \mathbf{s}, \mathbf{k}_2 \cdot \mathbf{s}, \ldots, \mathbf{k}_K \cdot \mathbf{s})$ (de Ruyter van Steveninck & Bialek, 1988). The number of relevant linear filters, $K$, is equal to the number of nonzero eigenvalues of the spike-triggered covariance matrix. A successful application of STC requires $P(\mathbf{s})$ to be Gaussian, and the number of filters $K$ be small (usually $\leq 3$) to ensure an adequate sampling of the filters and the nonlinearity $f$, given the data obtained in the typical experiment (however when inferring only the linear part of such models as many as 14 filters have been estimated (Rust et al., 2004)). STC has been used successfully, for example, to understand the computations performed by motion sensitive neurons in the blowfly (Bialek & de Ruyter van Steveninck, 2005), to map out the sensitivity to full-field flickering stimuli in salamander retinal ganglion cells (Fairhall et al., 2006), to explore contrast gain control (Rust et al., 2004; Schwartz et al., 2002), and to understand adaptation in the rodent barrel cortex (Maravall et al., 2007).

Before moving on, it seems appropriate to return once more to the Wiener formalism and contrast it with spike-triggered methods for recovering LN models. The underlying assumptions of the two approaches may seem substantially different: first, because of the presence of the nonlinear (N) transformation in the LN model, and second, because the output of the LN model is usually taken to predict the rate of a stochastic point process, while Wiener series is intended for analyzing deterministic systems (Wiener, 1958). Nevertheless, it is easy to see that when uncorrelated (i.e. *white*) Gaussian noise is used to extract the filters of the LN model using spike triggered average (STA) and spike triggered covariance (STC), STA and STC also provide unbiased estimates (up to a scaling factor) of first- and second-order Wiener kernels. The difference arises in subsequent analysis steps: in case of LN models, STA and STC are used solely as dimensionality reduction steps to identify the relevant subspace of the stimuli in which the nonlinear transformation acts, while in the Wiener formalism, STA and STC literally are the first two terms in a functional expansion that provides the best least-squares fit to the observed firing rate. Victor & Johannesma (1986) have further demonstrated that the Wiener formalism is a special case of a general probabilistic maximum entropy framework for describing joint distributions of stimuli and responses. In this framework, for example, the classic Wiener formalism is recovered if the stimulus distribution is Gaussian, and the response variable is also Gaussian with additive noise. If, on the other hand, the output variable is binary (spike / no-spike), the same maximum entropy approach reduces to identifying LN-type models with exponential nonlinearities.

While powerful and simple to use, spike-triggered covariance (STC) only works if Gaussian stimuli are employed, and is feasible only if $K$ is small. The Gaussian ensemble can be a serious restriction for neurons that do not respond well (or at all) to unstructured stimuli; furthermore, we are likely to miss several neural mechanisms that depend on naturalistic statistical structure, such as correlations, intermittency etc, if the neuron responds to Gaussian stimulation. A versatile method should therefore be able to successfully infer the multiple-filter dependence of a neuron probed with a stimulus of arbitrary complexity. Maximally informative dimensions (MID) (Sharpee et al., 2004) or likelihood inference for single-filter generalized linear models (Gerwinn et al., 2010; Paninski, 2004; Pillow, 2007; Pillow et al., 2008; Truccolo et al., 2004) have been used to this end when the dependence is linear, but the attempts to incorporate full quadratic stimulus dependence have been less common.

There are several instances of quadratic stimulus dependence. Let us consider a situation where the neuron has a vanishing spike-triggered average, as with a complex cell, non–phase–locked auditory neurons (Recio-Spinoso et al., 2005), or motion-sensitive neurons. In these cases a natural starting point would be a search for more than a single linear filter. For a model complex cell in the visual cortex, we would find two phase-shifted vectors $\mathbf{k}_1$ and $\mathbf{k}_2$ that together form a quadrature pair, such that the most informative variable concerning the neuron's firing is the "power,"

$$r(\mathbf{s}) = f\left[(\mathbf{k}_1 \cdot \mathbf{s})^2 + (\mathbf{k}_2 \cdot \mathbf{s})^2\right]. \tag{1}$$

Similarly, models of contrast gain control in the retina also include sensitivity to second-order features in the

**(a) Flickering variance stimulus**  **(b) Linear filters for different C**  **(c) Nonlinearities for different C**
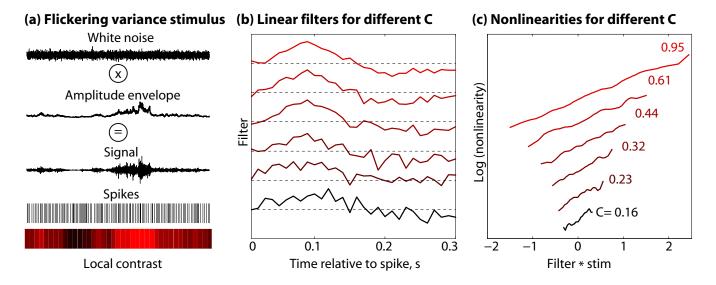


FIG. 2 **A synthetic contrast-adapting neuron probed with the "flickering variance" stimulus.** The instantaneous spiking rate is given by $r(t) = f(\mathbf{k}_0 \cdot \mathbf{s}(t) + \mathbf{s}(t)^T \mathbf{Q}\mathbf{s}(t) + \mu)$, where $f(\cdot) = \log(1 + \exp(\cdot))$, $\mu$ is an offset (bias), and the quadratic kernel $\mathbf{Q}$ is a rank 2 matrix with a quadrature eigenvector pair. **(a)** The stimulus is sampled at $\Delta = 1$ ms scale and is given by $s(t) = \exp(A(t))w(t)$, where $w(t)$ is given by uncorrelated white noise of fixed variance, and $A(t)$ is a gaussian noise process with correlation time $\tau_c = 1$ s. The stimulus can be chopped into segments of duration $\tau \leq \tau_c$, which can be sorted by local contrast $C$ (intensity of red). Spike-triggered average analysis can be applied to recover effective LN models for all stimulus segments sharing the same local contrast. **(b)** The linear filters recovered at various contrast levels (shade of red; filters displaced along vertical axis for easier readability). At lower contrasts the neuron produces less spikes, making the filter estimate more noisy, but the filter shape is constant across a range of $C$ and closely approximates the model filter $\mathbf{k}_0$. **(c)** The nonlinearities for different contrast levels $C$ (plot legend, shades of red; the nonlinearities displaced along vertical axis for easier readability). The slope of the nonlinearity decreases with increasing contrast (although the adaptation is not perfect, in this example), to prevent quick saturation of the response at high $C$.

stimulus, with the spiking probability of the form (Schwartz et al., 2002),

$$r(\mathbf{s}) = \frac{f(\mathbf{k_0} \cdot \mathbf{s})}{\sum_{i=1}^M w_i (\mathbf{k}_i \cdot \mathbf{s})^2 + \sigma^2}, \tag{2}$$

where the quadratic terms in the denominator scale down the gain at high contrast (in this case however, the neuron has a non-vanishing linear filter $\mathbf{k}_0$). A simulated model neuron showing contrast adaptation is shown in Fig. 2a, featuring both the first- and second-order stimulus sensitivity. The model neuron is probed with a "flickering variance" stimulus, in which the variance of the white noise (with a very short correlation time) is dynamically modulated by a noise process correlated across a longer timescale (c.f. Fairhall et al. (2001)). With this synthetic stimulus, the separation of timescales allows us to partition the stimulus into chunks with approximately constant variance in luminance, $\sigma_L^2$. This variance is directly related to the temporal contrast, $C = \sigma_L / \bar{L}$, because the average mean light intensity $\bar{L}$ is kept constant. Within each stimulus segment, we can use STA to recover the LN model, as shown in Figs. 2b,c. Our real goal, however, is to infer a joint model valid across the whole stimulus, and to do so ultimately with naturalistic stimuli with scale-free power spectra, where no clear separation exists between the fast fluctuation and slow variance modulation.

We can describe these and similar examples by a generic "quadratic" model neuron which is sensitive to a second-order function of the input (parametrized by a real, symmetric matrix $\mathbf{Q}$) in addition to the linear projection (parametrized by the filter $\mathbf{k}_0$):

$$r(\mathbf{s}) = f(\mathbf{k}_0 \cdot \mathbf{s}, \ \mathbf{s}^T \mathbf{Q}\mathbf{s}). \tag{3}$$

Graphically, while a threshold LN model with a linear filter corresponds to a classifier whose separating hyperplane is perpendicular to the filter, the proposed LN model with a threshold nonlinearity and a quadratic filter $\mathbf{Q}$ is selective for all stimuli that lie outside an $N$-dimensional ellipsoid whose axes correspond to the eigenvectors of $\mathbf{Q}$.

For the contrast gain control model described in Eq (2) the matrix $\mathbf{Q}$ is of rank $M$, with eigenvalues $w_i$ and eigenvectors $\mathbf{k}_i, i > 0$. The complex cell example described in Eq (1) has $\mathbf{k}_0 = 0$ and $\mathbf{Q} = \sum_{i=1}^2 \mathbf{k}_i \mathbf{k}_i^T$; in other

words, $\mathbf{Q}$ is a rank 2 matrix. While these examples feature quadratic dependences involving matrices of low rank, it is possible to extend these models to biologically relevant cases where the matrix does not have to be low rank (Rajan & Bialek, 2012). For example, the probability of spiking could be a nonlinear function of the "power" $p(t)$, $r(t) = f[p(t)]$, where the power is given by:

$$p(t) = \int d\tau f_2(\tau) \left[ \int dt' f_1(t - \tau - t')s(t') \right]^2;$$

(4)

here $s(t)$ is the stimulus, and $f_1$ and $f_2$ are linear filters, such as those used to describe non-phase-locked auditory neurons. If the smoothing time of the second filter $f_2$ is larger than that of the first filter $f_1$, it has been shown in (Rajan & Bialek, 2012) that the quadratic kernel $\mathbf{Q}$ for this model has a rich (full-rank) spectrum.

In the next section we review methods that permit inference of low- or full-rank quadratic kernels, $\mathbf{Q}$.

## III. INFERRING QUADRATIC STIMULUS DEPENDENCE FROM DATA

Every real, symmetric matrix can be spectrally decomposed into $\mathbf{Q} = \sum_{i=1}^{N} \lambda_i \mathbf{k}_i \mathbf{k}_i^{\mathrm{T}}$. The response of the quadratic model is thus $r = f \left[ \sum_{i=1}^{N} \lambda_i (\mathbf{k}_i \cdot \mathbf{s})^2 \right]$, explicitly demonstrating that quadratic models are special cases of the LNLN cascade, where the first linear stage involves applying the filters $\mathbf{k}_i$, the first nonlinear stage squares the projections, the second linear stage is a summation with weights $\lambda_i$, and the last nonlinear transformation is $f(\cdot)$. The spectral decomposition implies that we could try recovering the quadratic dependence of $\mathbf{Q}$ in Eq (3) by, for example, multidimensional MID (see Table I), hoping to infer all $\{\mathbf{k}_i\}$ as orthogonal informative dimensions. While formally true, this is infeasible in practice because maximizing the mutual information would involve sampling $N$-dimensional distributions from stimulus samples that are limited in number by the number of spikes (Sharpee et al., 2004). The same sampling problem would reappear when trying to estimate the nonlinearity, $f(\mathbf{k}_1 \cdot \mathbf{s}, \mathbf{k}_2 \cdot \mathbf{s}, \ldots, \mathbf{k}_N \cdot \mathbf{s})$.

To address this problem efficiently, we formulate the inference problem by explicitly assuming quadratic dependence on the stimulus: in this case, the stimulus immediately gets projected down to a single scalar variable $x = \mathbf{s}^{\mathrm{T}} \mathbf{Q} \mathbf{s}$, meaning that information-theoretic quantities, the likelihood, as well as the nonlinearity $f(\cdot)$ will only depend on the stimulus through $x$. This makes inference problem tractable even when $\mathbf{Q}$ is of high rank. Clearly, this advantage is gained by assuming that projections onto eigenvectors of $\mathbf{Q}$ combine as a sum of squares. This assumption is not a mere mathematical convenience: as we have shown previously, well-known phenomena of phase invariance, adaptation to local contrast or sensitivity to the signal envelope are all examples of true second-order stimulus sensitivity in real neurons. Additionally, response phenomena in the visual cortex grouped together as relating to the *non-classical receptive field* could also be manifestations of quadratic or higher-order sensitivity (Zetzsche & Nuding, 2005).

### A. Finding quadratic filters using information maximization

Despite their utility and simplicity, spike-triggered methods require the use of statistically simple stimuli and in particular, exclude the use of stimuli with naturalistic statistics, e.g. those with $1/f$ spectra, non-Gaussian histograms and/or high-order correlations. This is a big challenge when studying neurons beyond the sensory periphery that are responsible for extracting high-order structure, or neurons unresponsive to white noise presentations, for example those in the auditory pathway. To address this issue and recover the filter(s) in an unbiased manner with an arbitrary stimulus distribution, maximally informative dimensions (MID) (Kouh & Sharpee, 2009; Sharpee et al., 2004, 2006) have been developed and utilized to recover simple cell receptive fields, among other examples. MID looks for a linear filter $\mathbf{k}$ that maximizes the information between the presence/absence of a spike and the projection $x$ of the stimulus onto $\mathbf{k}$, $x = \mathbf{k} \cdot \mathbf{s}$. Information per spike is then given by the Kullback-Leibler divergence of $P(x|\text{spike})$, the *spike-triggered distribution* (the distribution of stimulus projections preceding the spike) and $P(x)$, the *prior distribution* (the overall distribution of projections):

$$I_{\text{spike}} = D_{KL} \left[ P(x|\text{spike}) || P(x) \right] = \int dx\, P(x|\text{spike}) \log_2 \frac{P(x|\text{spike})}{P(x)}.$$

(5)

Given the spike train and the stimulus, finding $\mathbf{k}$ becomes an information optimization problem in $I_{\text{spike}}$ that can be solved using various annealing methods, although the existence of local extrema makes this a nontrivial task.

Spike-triggered methods and MID do not explicitly assume a form for the nonlinearity $f(\cdot)$ in the LN model; instead, they provide unbiased estimates of the filter(s), and once the filters are known, the nonlinearity can be reconstructed

using the Bayes' rule from sampled spike-triggered and prior distributions:

$$f(x) \propto P(\text{spike}|x) = \frac{P(x|\text{spike})P(\text{spike})}{P(x)}, \tag{6}$$

where $P(\text{spike})$ is directly proportional to the average firing rate during the experiment.

In classical MID, one finds a (set of) linear filter(s) by maximizing Eq. (5) with respect to $\mathbf{k}$. In Rajan & Bialek (2012), this approach was extended to quadratic stimulus sensitivity, as follows. A quadratic filter $\mathbf{Q}$ can be reconstructed from an observed spike train by maximizing the information in Eq (5), where $x$ is now given by $x = \mathbf{s}^{\mathrm{T}}\mathbf{Q}\mathbf{s}$. Taking a derivative of Eq (5) with respect to $\mathbf{Q}$ gives us a gradient,

$$\nabla_{\mathbf{Q}} I = \int dx \, P_{\mathbf{Q}}(x) \left[ \langle \mathbf{s}\mathbf{s}^{\mathrm{T}} | x, \text{spike} \rangle - \langle \mathbf{s}\mathbf{s}^{\mathrm{T}} | x \rangle \right] \frac{d}{dx} \left( \frac{P_{\mathbf{Q}}(x|\text{spike})}{P_{\mathbf{Q}}(x)} \right), \tag{7}$$

where $\langle \cdot \rangle$ indicates averaging over the spike-triggered and prior distributions respectively, and the subscript $\mathbf{Q}$ makes the dependence of the probability distributions explicit. Only the symmetric part of $\mathbf{Q}$ contributes to $x$, and the overall scale of the matrix is irrelevant to the information, making the number of free parameters $N(N+1)/2 - 1$.

To learn the "Maximally Informative Stimulus Energy" or the quadratic filter $\mathbf{Q}$, we can ascend the gradient in successive learning steps (Rajan & Bialek, 2012),

$$\mathbf{Q} \to \mathbf{Q} + \gamma \, \nabla_Q I. \tag{8}$$

The probability distributions within the gradient are obtained by computing $x$ for all stimuli, choosing an appropriate binning for the variable $x$, and sampling binned versions of the spike-triggered and prior distributions. The $\langle \mathbf{s}\mathbf{s}^{\mathrm{T}} \rangle$ averages are computed separately for each bin; and the integral in Eqs (5,7) and the derivative in Eq (7) are approximated as a sum over bins and as a finite difference, respectively. To deal with local maxima in the objective function, we use a large starting value of $\gamma$ and gradually decrease $\gamma$ during learning. This basic algorithm can be extended by using kernel density estimation and stochastic gradient ascent/annealing methods, but we do not report these technical improvements here.

It is possible to select an approximate linear basis in which to expand the matrix $\mathbf{Q}$, by writing

$$\mathbf{Q} = \sum_{\mu=1}^{M} \alpha_\mu \mathbf{B}^{(\mu)}. \tag{9}$$

The basis can be chosen so that increasing the number of basis components $M$ allows the reconstruction of progressively finer features in $\mathbf{Q}$. We considered as $\{\mathbf{B}^{(\mu)}\}$ a family of Gaussian bumps that tile the space of the $N \times N$ matrix $\mathbf{Q}$ and whose scale (standard deviation) is inversely proportional to $\sqrt{M}$. For $M \to N^2/2$ the matrix set becomes a complete basis, allowing every $\mathbf{Q}$ to be exactly represented by the vector of coefficients $\alpha$. In such a matrix basis representation, the learning rule becomes

$$\alpha_\mu \to \alpha_\mu + \gamma \sum_{i,j=1}^{N} \frac{\partial I}{\partial \mathbf{Q}_{ij}} \mathbf{B}_{ij}^{(\mu)}, \tag{10}$$

where applying the chain rule on $\nabla_{\mathbf{Q}} I$ yields the $\text{Trace}[\nabla_{\mathbf{Q}}(\alpha) \cdot \mathbf{B}]$ update term at each step.

We illustrate this approach with two examples. In the first example we make use of the matrix basis expansion from Eq (9) to infer a quadratic kernel $\mathbf{K}$ that is of arbitrarily high rank. For $\mathbf{K}$ we used a highly-structured $500 \times 500$ matrix as shown in Fig. 3(a). While this is not an example of a receptive field from a real neuron, it illustrates the validity of the approach even when the response has an atypical and highly structured dependence on the stimulus. The stimuli were natural image clips from the Penn Natural Image database, flattened into a high-dimensional vector representation $\mathbf{s}$ (Tkačik et al., 2011), and the spikes were generated by thresholding the term $\mathbf{s}^{\mathrm{T}}\mathbf{K}\mathbf{s}$. Gaussian basis matrices, similar to the 225 shown in Fig. 3(b) were used to expand the quadratic kernel, reducing the number of optimization parameters from $\sim 2.5 \times 10^5$ to a few hundred. We start the gradient ascent with a large $\gamma$ value of 1 and progressively scale it down to 0.1 near the end of the algorithm; Fig. 3(e) shows the information plateauing in about 20 learning steps. The maximally informative quadratic filter reconstructed from 400 basis coefficients is shown in Fig. 3(d). Figure 3(c) demonstrates how the root-mean-squared reconstruction error systematically decreases as the number of basis functions $M$ is increased from 4 to 400, improving precision. Insets show 2 inferred matrices with $M = 100$ (corresponding to the first dot) showing a marked improvement with $M = 225$ (corresponding to the second red dot). Reconstruction error drops to $\sim 1\%$ for $M = 400$.

**(a) Model matrix, K**

**(b) Basis functions**

**(c) Reconstruction error**

100 bases

225 bases

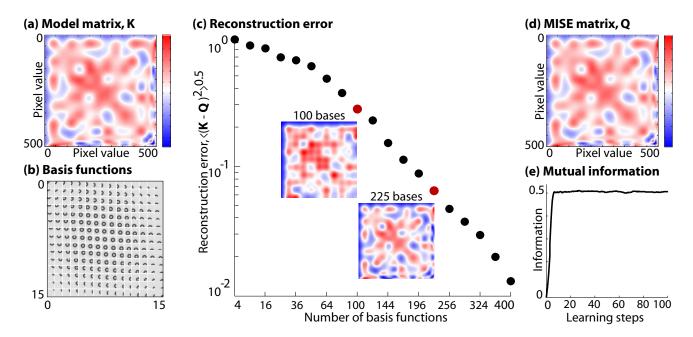**(d) MISE matrix, Q**

**(e) Mutual information**

FIG. 3 **Reconstructing a high rank quadratic filter using stimuli extracted from natural scenes.** **(a)** A complex high-rank randomly generated matrix $\mathbf{K}$ will be used as a quadratic filter of a model cell that fires whenever $\mathbf{s}^T\mathbf{K}\mathbf{s}$ exceeds a fixed threshold. $\mathbf{K}$ is thus the true quadratic filter for our threshold LN model neuron. **(b)** A collection of 225 Gaussian matrix basis functions whose peaks densely tile the matrix space; a trial matrix is constructed as a linear sum (with coefficients $\{\alpha_\mu\}$) of the basis matrices, and information optimization is performed over $\{\alpha_\mu\}$. **(c)** The normalized reconstruction error, shown in black dots, decreases as the number of basis functions $M$ increases from 4 to 400; with enough data perfect reconstruction is possible as $M$ approaches the number of independent pixels in $\mathbf{K}$. The two red dots show reconstructions with $M = 100$ or $M = 225$ basis functions, respectively. **(d)** The reconstructed, maximally informative matrix kernel $\mathbf{Q}$ after maximizing mutual information using 400 basis functions. **(e)** Mutual information increases as learning progresses in steps given by Eq. (8), peaks at step 40 and remains unchanged thereafter. Learning step 100 is the point where the maximally informative $\mathbf{Q}$ is extracted.

In contrast to standard MID where the number of spikes required grows exponentially in the number of filters extracted, the data requirement for this approach is proportional to the square of the stimulus dimension for a matrix kernel with no additional structural simplifications (these data requirement- and performance-related issues are explored in detail in Rajan & Bialek (2012)). For the examples shown in the paper, expansion in matrix basis reduces this number to the order of stimulus dimension, making this procedure pertinent for experimentalists.

The second example shows the MISE analysis of the synthetic neuron presented in Fig. 2 where stimulus-response relationship is more biologically realistic, through a smooth nonlinear function $f$ and both a linear as well as a quadratic kernel. The analysis is applied to the flickering variance stimulus without partitioning it into regions of fixed contrast. With $\sim 2 \times 10^4$ spikes, the method recovers the linear filter $\mathbf{k}_0$ as well as the quadratic kernel, which turns out to have the two dominant eigenvectors $\mathbf{k}_1, \mathbf{k}_2$, corresponding to the quadrature pair of filters used to construct $\mathbf{Q}$, as shown in Fig. 4b.

These examples show that quadratic filters can be extracted using information maximization for both low-rank and full-rank matrices, under natural stimulation and with a realistic numbers of spikes. Importantly, for cases where the stimulus sensitivity is both linear and quadratic, MISE does not explicitly assume that the effects of two filtering operations are additive, i.e. that $x = \mathbf{k}_0 \cdot \mathbf{s} + \mathbf{s}^T\mathbf{Q}\mathbf{s}$; rather, the dependence can be an arbitrary 2D nonlinear function, $f(\mathbf{k}_0 \cdot \mathbf{s}, \mathbf{s}^T\mathbf{Q}\mathbf{s})$. Unlike the quadratic generalizations of GLM presented below, this allows MISE to fully recover forms of contrast gain control that have a parametric form similar to Eq. (2).

## B. Finding quadratic filters using maximization of noise entropy

Another information-theoretic approach for inferring single neuron sensitivities is derived from the principle of noise entropy maximization (Fitzgerald et al., 2011a,b; Globerson et al., 2009). Suppose that the spiking or silence of a chosen neuron in a time bin indexed by $t$ is represented by a binary variable $y_t \in \{0, 1\}$. From data, we can reliably estimate certain statistics of the neural response, such as the average spiking rate $\langle y_t \rangle_t$, the spike-triggered average $\langle y_t \mathbf{s}(t) \rangle_t$, or the spike-triggered covariance $\langle y_t \mathbf{s}(t)\mathbf{s}(t)^T \rangle_t$, where the brackets $\langle \cdot \rangle_t$ denote averaging across the duration

of the experiment. In general, all these statistics are of the form $\langle O_\mu(\mathbf{s})y_t\rangle_t$, where $\mu$ indexes the different operators whose expectation values we are computing.

The crucial step is to look for maximum entropy approximations to $P(y|\mathbf{s})$, the distribution of the (binary) neural response given the stimulus. Maximum entropy distributions are as unstructured (random, therefore parsimonious) as possible with the constraint that they exactly reproduce the measured expectation values of a chosen set of statistics, $\{O_\mu\}$ (Jaynes, 1957a,b). When the variable $y$ is binary, it can easily be shown that these distributions have the form of the logistic function,

$$P(y = 1|\mathbf{s}) = \frac{1}{1 + e^{-F(\mathbf{s})}}, \tag{11}$$

where $F$ resembles the free energy in statistical physics:

$$F(\mathbf{s}) = \sum_\mu \lambda_\mu O_\mu(\mathbf{s}), \tag{12}$$

and $\lambda_\mu$ are the Lagrange multipliers that have to be set such that the set of statistics measured in the data equals the expectation values of the same operators under distribution $P$, i.e. $\langle O_\mu(\mathbf{s})y\rangle_P = \langle O_\mu(\mathbf{s})y\rangle_t$. To apply this general framework to the inference of quadratic filters, the authors of Fitzgerald et al. (2011b) choose the mean firing rate, STA and STC as constraints, which yields the following response distribution:

$$P(y = 1|\mathbf{s}) = \frac{1}{1 + \exp(\mu + \mathbf{k}_0 \cdot \mathbf{s} + \mathbf{s}^{\mathrm{T}}\mathbf{Q}\mathbf{s})}, \tag{13}$$

where $\{\mu, \mathbf{k}_0, \mathbf{Q}\}$ act as the Lagrange multipliers $\lambda_\mu$ conjugated to the operators $\{y, y\mathbf{s}, y\mathbf{s}\mathbf{s}^{\mathrm{T}}\}$. Numerically, the task is to solve for parameters $\{\mu, \mathbf{k}_0, \mathbf{Q}\}$ that satisfy a set of constraints: $\langle y\rangle_t = \langle y\rangle_P$ (matching the measured mean firing rate to that of the model), $\langle y\mathbf{s}\rangle_t = \langle y\mathbf{s}\rangle_P$ (matching the measured STA to that of the model), and $\langle y\mathbf{s}\mathbf{s}^{\mathrm{T}}\rangle_t = \langle y\mathbf{s}\mathbf{s}^{\mathrm{T}}\rangle_P$ (matching the measured STC to that of the model). This is a convex optimization task and can be solved by conjugate gradient descent.

An attractive feature of this approach emerges when we rewrite the information per spike $I(\text{spike}; \mathbf{s})$ as a difference between the total and the noise entropy as follows:

$$I(\text{spike}; \mathbf{s}) = \sum_\mathbf{s} P(\mathbf{s}) \sum_y P(y|\mathbf{s}) \log_2 \frac{P(y|\mathbf{s})}{P(y)} = S[P(y)] - \langle S[P(y|\mathbf{s})]\rangle_\mathbf{s}, \tag{14}$$

where $S[P(x)] = -\sum_x P(x) \log_2 P(x)$ is the entropy of $P(x)$. The first term (total entropy) is fully determined by the mean spiking rate $\langle y\rangle_t$, $S[P(y)] = -\langle y\rangle_t \log_2\langle y\rangle_t - (1 - \langle y\rangle_t) \log_2(1 - \langle y\rangle_t)$ because $y$ is a binary variable. The mean firing rate is one of the statistics constrained in the model for $P(y|\mathbf{s})$, ensuring consistency. Since our model for $P(y|\mathbf{s})$ has maximum entropy given the observed constraints, we are effectively setting an upper bound on the noise entropy $\langle S[P(y|\mathbf{s})]\rangle_\mathbf{s}$, and therefore a lower bound on the mutual information $I$. As more and more statistics $O(\mathbf{s})$ are included as constraints into the maximum entropy model for Eq. (11), the noise entropy must progressively drop and information increase towards the true value (which is bounded by the output entropy). At the point where this lower bound on information meets the actual information per spike (which can be empirically estimated from, e.g., repeated stimulation (Brenner et al., 2000)), we obtain the complete list of the relevant stimulus statistics $\{O_\mu\}$ that characterize the sensitivity of the neuron.

In Fitzgerald et al. (2011b), the authors show that this framework is applicable for inferring quadratic neural filters on synthetic and real data, and compare it to MID. This method is applicable to any stimulus ensemble, but requires assumptions beyond those needed for MID or MISE: namely, that the nonlinear function is logistic, and that the contributions of the linear and quadratic filters add. The advantage of the method is that the problem is convex, does not suffer from the exponential curse of dimensionality (like multi-dimensional MID), and is flexible, permitting various new constraints (beyond the STA and STC) to be used in constructing models for the stimulus-conditional distribution $P(y|\mathbf{s})$.

## C. Finding quadratic filters in a likelihood framework: GLM extensions and Bayesian STC

A powerful technique for modeling neural spiking behavior is the generalized linear model (GLM) framework (Paninski, 2004; Truccolo et al., 2004). Recently GLM has been used to account for the stimulus sensitivity, dependence on spiking history, and connectivity in a population of 27 retinal ganglion cells in the macaque retina (Pillow et al.,

2008). For a single neuron, the model assumes that the instantaneous spiking rate $r(t)$ is a nonlinear function $f$ of a sum of contributions,

$$r(t) = f\left[\mathbf{k} \cdot \mathbf{s}(t) + \mathbf{q} \cdot \mathbf{y}(t_-) + \mu\right], \tag{15}$$

where $\mathbf{k}$ is a linear filter acting on the stimulus $\mathbf{s}$, $\mathbf{q}$ is a linear filter acting on the spiking history $\mathbf{y}(t_-)$ of the neuron, and $\mu$ is an offset or an intrinsic bias towards spiking or silence. When the stimulus and the spike train are discretized into time bins of duration $\Delta$, the probability of observing (an integral number of) $y_t$ spikes is Poisson, with a mean given by $r_t \Delta$ (where the subscript indexes the time bin). Here, we neglect the history dependence of the spikes (with no loss of generality) and focus instead on the stimulus dependence; since each time bin is conditionally independent given the stimulus (and past spiking), the log likelihood for any spike train $\{y_t\}$ is (Pillow, 2007):

$$\log P(\{y_t\}|\mathbf{s}) = \sum_t y_t \log r_t - \Delta \sum_t r_t + c, \tag{16}$$

where $c$ is independent of both $\mu$ and $\mathbf{k}$. This likelihood can be maximized with respect to $\mu$ and $\mathbf{k}$ (and optionally, with respect to $\mathbf{g}$) given adequate number of spikes, providing an estimate of the filters from neural responses to complex, even natural stimuli. In contrast to maximally informative approaches, such as the stimulus energy derived in Section III.A (Rajan & Bialek, 2012), the functional form of the nonlinearity $f$ is an explicit assumption in likelihood-based methods like GLM. For specific classes of the function $f$, such as $f(z) = \log[1 + \exp(z)]$, $\exp(z)$ or $[1 + \exp(z)]^{-1}$, the likelihood optimization problem is convex and gradient ascent is guaranteed to find a unique global maximum.

While the tractability consequent to convexity of the objective function is a big strength of this approach, the disadvantage is that if the chosen nonlinearity $f$ is different from the true function $f'$ used by the neuron, the filters inferred by maximizing likelihood in Eq (16) could be biased. If we relax the stringent requirement for convexity, we can choose more general nonlinear functions for the model, for example by parametrizing the nonlinearity in a point-wise fashion and inferring it jointly with the filters. For this discussion however, we assume that $f$ has been selected from the specific class of nonlinearities guaranteed to yield a convex likelihood function.

How can we extend GLM to situations where the neuron's response is more complex than a single linear projection of the stimulus? We will start with a proposal and follow up with a closely related formulation of Park & Pillow (2011) developed in parallel, which has provided a more complete analysis and several interesting extensions. One possibility is to expand the stimulus clip $\mathbf{s}$ of dimension $N$ into a larger space first, for instance by forming $\mathbf{s}\mathbf{s}^{\mathrm{T}}$ (of dimension $N \times N$), and then operate on this object with a filter, i.e., $\sum_{i,j=1}^{N}(s_i s_j)Q_{ij}$. Such a term can be added to the argument of $f$ in the model exemplified in Eq (15). Specifically, we propose a "Generalized Quadratic Model" of the following form,

$$r(t) = f\left[\mathbf{k} \cdot \mathbf{s}(t) + \mathbf{s}^{\mathrm{T}}(t)\mathbf{Q}\mathbf{s}(t) + \mathbf{g} \cdot \mathbf{y}(t_-) + \mu\right]. \tag{17}$$

If we want to retain convexity, we cannot expand $\mathbf{Q}$ in its eigenbasis and infer its vectors by maximizing the likelihood directly, because the eigenvectors appear quadratically. However, we can expand $\mathbf{Q}$ into a weighted sum of matrix basis functions, as in Eq (9), making the argument of $f$ a linear function of basis coefficients $\alpha_\mu$,

$$r(t) = f\left(\mathbf{k} \cdot \mathbf{s}(t) + \sum_{\mu=1}^{M} \left[\mathbf{s}^{\mathrm{T}}(t)\mathbf{B}^{(\mu)}\mathbf{s}(t)\right]\alpha_\mu + \mathbf{g} \cdot \mathbf{y}(t_-) + \mu\right). \tag{18}$$

Existing methods for inferring GLM parameters (Pillow et al., 2008) can be used to learn both the linear filter and the quadratic filter $\mathbf{Q}$ efficiently. After extracting $\mathbf{Q}$ we can check if a few principal components account for most of its structure (this is equivalent to checking whether $\mathbf{Q}$ is indeed a low rank matrix). In sum, this procedure provides a way of extracting multiple filters with GLM that is analogous to diagonalizing the spike-triggered covariance matrix on the Gaussian stimulus ensemble.

We have implemented such a quadratic extension to the GLM and applied it to the flickering variance stimulus shown in Fig. 2. The results are shown in Fig. 4a. The quadratic kernel correctly recovers a quadrature pair of filters; we similarly recover the correct linear filter $\mathbf{k}_0$. While this method is restricted to a linear combination of first- and second-order filters within the nonlinearity, the distinct advantage over MISE is that the inference problem is convex with the appropriate nonlinearity.

Park & Pillow (2011) consider an exponentiated general quadratic function of the following form (rewritten in the notation of this paper):

$$r(\mathbf{s}) = \exp\left(\mathbf{s}^{\mathrm{T}}\mathbf{Q}\mathbf{s} + \mathbf{k}_0 \cdot \mathbf{s} + \mu\right). \tag{19}$$
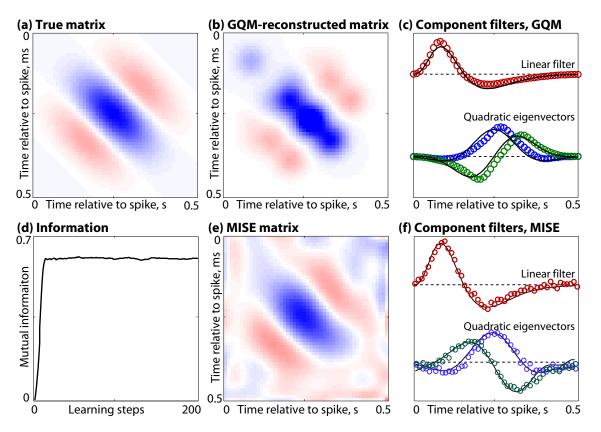
FIG. 4 **Recovering the synthetic model of the contrast gain control cell using the flickering variance stimulus.**
The spikes were simulated using the model presented in Fig. 2. **(a)** The true quadratic kernel, $\mathbf{Q}$, of the model, is a matrix
of rank 2 with the two filters combining into quadrature to estimate the signal "power" or variance. **(b)** The reconstructed
kernel using the quadratic extension of the GLM; the space of matrices was spanned by a 85-dimensional basis of Gaussian
bumps (some of the granularity can still be seen in the reconstruction). The dominant eigenvectors of the inferred matrix are
shown in **(c)** in blue and green (solid black lines show the true values); shown is also the recovered linear filter (red circles)
and its true value (solid black line). The inference of the same model using MISE shows quick convergence in **(d)** and the
recovered quadratic kernel in **(e)**. **(f)** The linear filter and the eigenvectors of the quadratic kernel recovered with MISE
(circles), compared to the true values (black solid line). Note that quadratic filter eigenvectors are only determined up to a
sign.

First, the authors show that under a Gaussian stimulus ensemble, the expected log likelihood can be expressed in terms
of the STA, STC, and the covariance matrix of the stimulus, and derive the closed–form expressions for maximum
likelihood estimates of the quadratic kernel, linear filter, and the bias. Next, the generalization to arbitrary stimuli
is achieved by numerically optimizing the true (as opposed to "expected") likelihood. In contrast to our suggestion
of using the matrix basis expansion (which becomes an implicit regularizer upon choosing the dimensionality of
the basis), Park and Pillow implement Bayesian regularization by imposing a prior on the quadratic kernel. This
important suggestion is implemented as follows.

The matrix is first decomposed into the eigensystem, $\mathbf{Q} = \sum_{i=1}^{N} \sigma_i \mathbf{w}_i \mathbf{w}_i^{\mathrm{T}}$, where $\mathbf{w}$ are not forced to have an $L_2$
norm of 1, and $\sigma_i = \pm 1$ to indicate whether the filter $i$ is excitatory or suppresive (as in STC (Schwartz et al.,
2006)). Then, a zero-mean Gaussian prior $\mathcal{N}(0, \alpha_i^{-1}I)$ is put on each eigenvector $\mathbf{w}_i$, where the hyperparameter $\alpha_i$
determines the variance of the elements of eigenvector $i$; $\alpha_i \to \infty$ corresponds to eliminating the direction $i$ from
the quadratic kernel and reducing its rank by 1. Next, an iterative algorithm is described for alternating between
optimizing the likelihood with respect to model parameters, and optimizing the evidence given the parameters with
respect to hyperparameters $\alpha_i$. This procedure correctly and efficiently identifies the rank of the quadratic kernels in
synthetic examples, providing an automatic alternative for distinguishing "significant" from sampling-noise-induced
eigenvectors in the STC and quadratic kernel inference. Finally, the authors show that Eq. (19) can be further
generalized at no additional computational cost from the exponentiated quadratic function to a wider class of elliptic
nonlinearities.

To summarize, the reviewed work shows that the Bayesian generalization of STC and the generalization of GLMs

to quadratic stimulus dependence yield equal probabilistic models for neural encoding that can be efficiently inferred for a restricted class of nonlinear functions. Attention needs to paid, however, to maintain the convexity of the procedure and deal with the large number of parameters in the quadratic kernel. To this end, basis expansions as well as regularization with Bayesian priors seem like feasible candidates.

## IV. DISCUSSION

While powerful conceptually, the notion that neurons respond to multiple projections of the stimulus onto orthogonal filters is difficult to turn into a tractable inference procedure when the number of filters is larger than a few. To address this concern, alternative encoding models have recently been proposed where the neuron can be sensitive to higher-order features in the stimulus. Instead of being described by multiple linear filters, the neuron's sensitivity is described by a single quadratic filter (and optionally an additional linear filter). We have reviewed several inference methods for such quadratic stimulus dependence: two based on information maximization and the other based on maximizing the likelihood in an extension of generalized linear models. With MISE, no assumptions are made about how the projection onto the quadratic filter combines with the linear filter projection, and how both map into the probability of spiking. This approach yields unbiased filter estimates under any stimulus ensemble, but requires optimization in a possibly rugged information landscape. Noise entropy maximization is a flexible, maximum-entropy based framework for modeling the probability of spiking given stimulus. It is computationally tractable and provides a convenient bound on the information per spike, but assumes a particular form of the nonlinearity. Alternatively, with a specific choice of nonlinearity and filter basis, likelihood inference within the GLM class can be extended to quadratic stimulus dependence while retaining the convexity of the objective function. By formulating the problem as Bayesian inference and choosing sparsifying priors for the quadratic filter, the true rank of the quadratic filter can also be inferred from data.

All these approaches for inferring quadratic stimulus dependence are complementary; as we show in the appendix, maximum likelihood and information maximization inference also provide consistent filter estimates under defined conditions. A possible way to benefit from the tractability of likelihood formulations and maximization of noise entropy could be to use them to initialize a more general search using information maximization, in the hope that this would avoid the problems with the rugged information landscape, and remove the restrictions on the additive combination of linear and quadratic features.

Examples of recent work establishing connections between higher-order structure of natural scenes and neural mechanisms beyond the sensory periphery (e.g. Karklin & Lewicki (2009); Tkačik et al. (2010)) make the development of corresponding methods for neural characterization, such as the ones presented here, very timely. Phenomena like phase invariance, adaptation to local contrast or sensitivity to signal envelope are widespread features of sensory neuron responses (Baccus & Meister, 2004; Hubel & Wiesel, 1965; Touryan et al., 2002). Moreover, as our abilities to record in vivo from the sensory systems of awake and behaving animals expand, so should the methods to analyze such recordings, where the relevant stimuli may no longer be perfectly controllable because of the animal's interaction with the environment (Kerr & Nimmerjahn, 2012). The methods presented here will help us systematically elucidate sensitivity to higher-order statistical features from responses of sensory neurons to natural stimuli.

### Appendix: The relationship between information theoretic and likelihood-based inference

We now demonstrate analytically that under rather general assumptions, the linear or quadratic filters obtained by maximizing mutual information match the filters inferred by maximizing the likelihood. We extend a reasoning we used previously in the context of inferring protein-DNA sequence-specific interactions in Kinney et al. (2007), to neural responses. See also Kouh & Sharpee (2009) and references therein for a similar demonstration.

In the following, $x$ remains the projection of the stimulus $\mathbf{s}$ onto the linear ($x_t = \mathbf{k} \cdot \mathbf{s}_t$) or quadratic ($x_t = \mathbf{s}_t^\mathrm{T} \mathbf{Q} \mathbf{s}_t$) filter, with time discretized in bins of duration $\Delta$ and indexed by subscript $t$. We collect all parameters that determine the filter into a vector $\theta_1$. Given a single $x_t$, $y_t$ spikes are generated according to a conditional probability distribution $\pi(y_t|x_t)$. This probability distribution is typically assumed to be Poisson with mean given by $f(x_t)$ in the case of GLM, but we take a different approach. We discretize $x_t$ into $x = 1, \ldots, K$ bins and parameterize $\pi(y_t|x_t)$, which is a $Y_{\max} \times K$ matrix, by a set of parameters $\theta_2$. Apart from assuming a cutoff value for the number of spikes per bin $Y_{\max}$ (which can always be chosen large enough to assign an arbitrarily low probability to observing $> Y_{\max}$ spikes in any real dataset) and a particular discretization of the projection variable $x$, we leave the probabilistic relationship $\pi(y|x)$ between the projection and spike count completely unconstrained. The transformation from the stimulus to

the spikes is then a Markov chain, fully specified by $\theta = \{\theta_1, \theta_2\}$,

$$\mathbf{s}_t \xrightarrow[\mathbf{k} \text{ or } \mathbf{Q}]{\theta_1} x_t \xrightarrow[\pi]{\theta_2} y_t. \tag{20}$$

The likelihood of the spike train $\{y_t\}$ given the stimulus $\mathbf{s}$ is $P(\{y_t\}|\mathbf{s}) = \prod_{t=1}^{T} \pi(y_t|x_t)$, where $T$ is the total number of time bins in the dataset. With $x$ discretized into $K$ bins, any dataset can be summarized by the count matrix $c_{yx} = \sum_{t=1}^{T} \delta(y, y_t)\delta(x, x_t)$, where $\delta$ is the Kronecker delta; note that $c_{yx} = T\tilde{p}(y, x)$, where $\tilde{p}$ is simply the empirical distribution in the data of observing $y$ spikes jointly with projection $x$. In terms of $c$, the likelihood of the observed spike train is $P(\{y_t\}|\mathbf{s}) = \prod_{y=0}^{Y_{\max}} \prod_{x=1}^{K} \pi(y|x)^{c_{yx}}$. Assuming that $x$ is adequately discretized and that $\pi$ is Poisson with mean $f(x)$, we will recover the generalized likelihood of Eq (16).

Suppose that we are only interested in inferring the filter (parametrized by $\theta_1$), but not the filter-to-spike mapping $\pi$ (parameterized by $\theta_2$). While avoiding any assumptions about the structure of $\pi$, we can integrate the likelihood over $\theta_2$ with some prior $P_p(\theta_2)$ such that

$$P(\{y_t\}|\mathbf{s}) = \int d\theta_2 \, P_p(\theta_2) \prod_{y,x} \pi(y|x)^{c_{yx}}. \tag{21}$$

This resulting likelihood, called the *model averaged likelihood*, is now only a function of $\theta_1$. The prior $P_p(\theta_2)$ can take many forms, but since we discretized $x$, thereby making $\pi(y|x)$ into a (conditional probability) matrix, the simplest choice for the prior is the *uniform prior*. In this case we set $\theta_2$ equal to the entries in $\pi(y|x)$ matrix and choose $P(\theta_2)$ to be uniform over all valid matrices $\pi$, such that the matrix entries are positive and the normalization constraint, $\sum_x \pi(y|x) = 1$ for every $x$, is enforced.

For any choice of priors we can rewrite Eq (21) as

$$P(\{y_t\}|\mathbf{s}) = \int d\theta_2 \, P_p(\theta_2) \exp\left[T \sum_{y,x} \tilde{p}(y, x) \log \pi(y|x)\right], \tag{22}$$

which can be reorganized into

$$P(\{y_t\}|\mathbf{s}) = \int d\theta_2 P_p(\theta_2) \exp\left[T\left\{\tilde{I}(y; x) - \tilde{S}(y) - \langle D_{KL}(\tilde{p}(y|x) \,||\, \pi(y|x))\rangle_{\tilde{p}(x)}\right\}\right]. \tag{23}$$

Here $\tilde{I}(y; x) = \sum_{y,x} \tilde{p}(y, x) \log \frac{\tilde{p}(y,x)}{\tilde{p}(y)\tilde{p}(x)}$ is the empirical mutual information between spike counts $y$ and the projection $x$, $\tilde{S}(y)$ is the empirical spike count entropy, and the "correction" term in brackets measures the average Kullback-Leibler divergence ($D_{KL}$) between the empirical and model conditional distributions. Importantly, only this correction term is a function of the $\pi$ and thus of $\theta_2$, and is affected by the prior $P_p(\theta_2)$ which is being integrated over; the other terms can be pulled outside of the integral. We can therefore write the per time bin log likelihood as

$$\mathcal{L} = \frac{1}{T} \log P(\{y_t\}|\mathbf{s}) = \tilde{I}(y; x) - \tilde{S}(y) - \Lambda, \tag{24}$$

where the correction is

$$\Lambda = -\frac{1}{T} \log \int d\theta_2 \, P_p(\theta_2) e^{-T\langle D_{KL}(\tilde{p}(y|x) \,||\, \pi(y|x))\rangle_{\tilde{p}(x)}}. \tag{25}$$

It is necessary to show that as the amount of data $T$ grows, the correction $\Lambda$ decreases for a given choice of prior distribution $P_p(\theta_2)$, and for the choice of uniform prior this is analytically tractable (Kinney et al., 2007). Intuitively, it is clear that as $T \to \infty$, the empirical distribution $\tilde{p}(y|x)$ converges to the true underlying distribution $p(y|x)$, and the integral becomes dominated by the extremal point $\theta_2^*$, such that, in the saddle point approximation,

$$\Lambda(T \to \infty) \sim \langle D_{KL}(p(y|x) \,||\, \pi^*(y|x))\rangle_{p(x)}. \tag{26}$$

The distribution $\pi^*(y|x)$ is the closest distribution to $p(y|x)$ in the space over which the prior $P_p(\theta_2)$ is nonzero. As long as the prior assigns a non-zero probability to any (normalized) distribution, the divergence in $\Lambda$ will decrease and $\Lambda$ will vanish as $T$ grows. The case in which $\Lambda$ does not decay occurs when the prior completely excludes certain distributions by assigning zero probability, while the data $p(y|x)$ precisely favors those excluded distributions.

Returning to the per time bin log likelihood $\mathcal{L}$ in Eq (24), as we decrease the time bin $\Delta$, we enter a regime where there is only 0 or 1 spike per bin, i.e., $y \in \{0, 1\}$. Then the empirical information per time bin $\tilde{I}(y; x)$ can be written as,

$$\tilde{I}(y; x) = \tilde{p}(y = 0)D_{KL}\left(\tilde{p}(x|y=0)||\tilde{p}(x)\right) + \tilde{p}(y = 1)D_{KL}\left(\tilde{p}(x|y=1)||\tilde{p}(x)\right), \tag{27}$$

that is,

$$\tilde{I}(y; x) = \tilde{p}(\text{silence})\tilde{I}_{\text{silence}} + \tilde{p}(\text{spike})\tilde{I}_{\text{spike}}. \tag{28}$$

If the information in the spike train is dominated by the information carried in spikes (Brenner et al., 2000), then the likelihood from Eq (24) becomes

$$\mathcal{L} = \tilde{p}(\text{spike})\tilde{I}_{\text{spike}} + \ldots, \tag{29}$$

where $\ldots$ are terms that either do not depend of the filter parameters $\theta_1$ (i.e. entropy of the spike counts $\tilde{S}(y)$), or vanish as the size of dataset grows ($\Lambda$).

The identity in Eq (29) is the sought-after connection between the inference using information maximization and the likelihood-based approach. In the limit of small time-bins, maximizing the information per spike $I_{\text{spike}}$ (in maximally informative approaches, as in (Sharpee et al., 2004) and Section III.A of this paper), on right-hand side of the identity, is the same as maximizing the *model averaged likelihood* $\mathcal{L}$ of Eq (24), on the left-hand side of the identity.

## Acknowledgments

## References

B Agüera y Arcas, AL Fairhall, & W Bialek (2003) Computation in a single neuron: Hodgkin and Huxley revisited. *Neural Comput* **15:** 1715–49.

B Agüera y Arcas & AL Fairhall (2003) What causes a neuron to spike? *Neural Comput* **15:** 1789–807.

SA Baccus & M Meister (2002) Fast and slow contrast adaptation in retinal circuitry. *Neuron* **36:** 909–919.

SA Baccus & M Meister (2004) Retina versus cortex; contrast adaptation in parallel visual pathways. *Neuron* **42:** 5–7.

MJ Berry 2nd & M Meister (1998) Refractoriness and neural precision. *J Neurosci* **18:** 2200–2211.

MJ Berry 2nd, IH Brivanlou, TA Jordan & M Meister (1999) Anticipation of moving stimuli by the retina. *Nature* **398:** 334–8.

W Bialek & RR de Ruyter van Steveninck (2005) Features and dimensions: Motion estimation in fly vision. *arxiv.org:q-bio/0505003.*

E de Boer & P Kuyper (1968) Triggered correlation. *IEEE Trans Biomed Eng* **15:** 169–179.

D Bölinger & T Gollisch (2012) Closed-loop measurements of iso-response stimuli reveal dynamic nonlinear stimulus integration in the retina. *Neuron* **73:** 333–346.

A Borst & M Egelhaaf (1987) Temporal modulation of luminance adapts time constant of fly movement detectors. *Biol Cybern* **56:** 209–215.

N Brenner, RR de Ruyter van Steveninck & W Bialek (2000) Adaptive rescaling maximizes information transmission. *Neuron* **26:** 695–702.

P Dayan & LF Abbott (2001) Theoretical neuroscience. MIT Press, Cambridge, MA.

JB Demb, K Zaghloul, L Haarsma & P Sterling (2001) Bipolar cells contribute to nonlinear spatial summation in the brisk-transient (Y) ganglion cell in mammalian retina. *J Neurosci* **21:** 7447–7454.

AL Fairhall, GD Lewen, W Bialek & RR de Ruyter van Steveninck (2001) Efficiency and ambiguity in an adaptive neural code. *Nature* **412:** 787–92.

AL Fairhall, CA Burlingame, R Narasimhan, RA Harris, JL Puchalla & MJ Berry 2nd (2006) Selectivity for multiple stimulus features in retinal ganglion cells. *J Neurophysiol* **96:** 2724–38.

JD Fitzgerald, LC Sincich & TO Sharpee (2011) Minimal models of multidimensional computations. *PLoS Comput Biol* **7:** e1001111.

JD Fitzgerald, RJ Rowekamp, LC Sincich & TO Sharpee (2011) Second order dimensionality reduction using minimum and maximum mutual information models. *PLoS Comput Biol* **7:** e1002249.

WS Geisler (2008) Visual perception and the statistical properties of natural scenes. *Annu Rev Psychol* **59:** 167–92.

W Gerstner & W Kistler (2002) Spiking neuron models: Single neurons, populations, plasticity. Cambridge, Cambridge University Press.

S Gerwinn, J Macke & M Bethge (2010) Bayesian inference for generalized linear models for spiking neurons. *Frontiers in Comput Neurosci* **4:** 12.

A Globerson, E Stark, E Vaadia & N Tishby (2009) The minimum information principle and its application to neural code analysis. *Proc Nat'l Acad Sci USA* **106:** 3490–3495.

E Granot-Atedgi, G Tkačik, R Segev & E Schneidman (2012) Stimulus-dependent maximum entropy models of neural population codes. *arXiv.org:1205.6438.*

T Gollisch & AVM Herz (2005) Disentangling sub-millisecond process within an auditory transduction chain. *PLoS Biol* **3:** e8.

T Gollisch & M Meister (2010) Eye smarter than scientists believed: neural computations in circuits of the retina. *Neuron* **65:** 150–64.

JH van Hateren (1992) Theoretical predictions of spatiotemporal receptive felds of fly LMCs, and experimental validation. *J Comp Physiol A* **171:** 157–170.

HK Hartline (1940) The receptive fields of optic nerve fibers. *Am J Physiol* **130:** 690–699.

A Hodgkin & A Huxley (1952) A quantitative description of membrane current and its application to conduction and excitation in nerve. *J Physiol* **117:** 500–544.

S Hong, B Agüera y Arcas & AL Fairhall (2007) Single neuron computation: from dynamical system to feature detector. *Neural Comput* **19:** 3133–72.

DH Hubel & TH Wiesel (1965) Receptive fields and functional architecture in two nonstriate visual areas (18 and 19) of the cat. *J Physiol* **28:** 229–289.

ET Jaynes (1957) Information theory and statistical mechanics. *Phys Rev* **106:** 620–630.

ET Jaynes (1957) Information theory and statistical mechanics II. *Phys Rev* **108:** 171–190.

Y Karklin & MS Lewicki (2009) Emergence of complex cell properties by learning to generalize in natural scenes. *Nature* **457:** 83–6.

J Keat, P Reinagel, R Clay Reid & M Meister (2001) Predicting every spike: a model for the responses of visual neurons. *Neuron* **30:** 803–17.

JND Kerr & A Nimmerjahn (2012) Functional imaging in freely moving animals. *Curr Op Neurobiol* **22:** 45–53.

JB Kinney, G Tkačik & CG Callan Jr (2007) Precise physical models of protein-DNA interaction from high-throughput data. *Proc Nat'l Acad Sci USA* **104:** 501–506.

C Koch (1999) *Biophysics of Computation: Information processing in single neurons.* Oxford University Press (New York).

M Kouh & TO Sharpee (2009), Estimating linear-nonlinear models using Renyi divergences. *Network* **20**: 49–68.

NA Lesica, T Ishii, GB Stanley & T Hosoya (2008) Estimating receptive fields from responses to natural stimuli with asymmetric intensity distributions. *PLoS ONE* **3:** e3060.

BN Lundstrom, S Hong, AL Fairhall (2008) Two computational regimes of a single-compartment neuron separated by a planar boundary in conductance space. *Neural Comput* **20:** 1239–60.

M Maravall, RS Petersen, AL Fairhall, E Arabzadeh & ME Diamond (2007) Shifts in coding properties and maintainance of information transmission during adaptation in barrel cortex. *PLoS Biol* **5:** e19.

PZ Marmarelis and VZ Marmarelis (1978) Analysis of physiological systems: The white-noise approach. New York: Plenum.

BP Olveczky, SA Baccus & M Meister (2007) Retinal adaptation to object motion. *Neuron* **56:** 689–700.

S Ostojic & N Brunel (2011) From spiking neuron models to linear-nonlinear models. *PLoS Comput Biol* **7:** e1001056.

Y Ozuysal & SA Baccus (2012) Linking the computational structure of variance adaptation to biophysical mechanisms. *Neuron* **73:** 1002-1015.

L Paninski (2003) Convergence properties of some spike-triggered analysis techniques. *Network* **14:** 437–464.

L Paninski, JW Pillow & EP Simoncelli (2004) Maximum likelihood estimation of a stochastic integrate-and-fire neural encoding model. *Neural Comput* **16:** 2533-2561.

L Paninski (2004) Maximum likelihood estimation of cascade point-process neural encoding models. *Network Comp Neural Syst* **15:** 243–62.

I Park & JW Pillow (2011), Bayesian spike-triggered covariance. *Advances in Neural Information Processing Systems (NIPS)* **24**:1692–1700 .

JW Pillow (2007) Likelihood-based approaches to modeling the neural code. In *Bayesian Brain: Probabilistic Approaches to Neural Coding,* eds K Doya, S Ishii, A Pouget & R Rao, pg. 53–70. MIT Press.

JW Pillow, J Shlens, L Paninski, A Sher, AM Litke, EJ Chichilnisky & EP Simoncelli (2008) Spatio-temporal correlations and visual signalling in a complete neural population. *Nature* **454:** 995–9.

JW Pillow & EP Simoncelli (2006) Dimensionality reduction in neural models: an information-theoretic generalization of spike-triggered average and covariance analysis. *J Vis* **6:** 414–428.

K Rajan & W Bialek (2012), Maximally informative "stimulus energies" in the analysis of neural responses to natural signals. *arXiv.org:*1201.0321.

J Rapela, G Felsen, J Touryan, JM Mendel & NM Grzywacz (2010) ePPR: a new strategy for the characterization of sensory cells from input/output data. *Network* **21:** 35–90.

A Recio-Spinoso, AN Temchin, P van Dijk, YH Fan & MA Ruggero (2005) Wiener-kernel analysis of responses to noise of Chinchilla auditory-nerve fibers. *J Neurophys* **93:** 3615–34.

RC Reid, JD Victor & RM Shapley (1997) The use of m-sequences in the analysis of visual neurons: linear receptive field properties. *Vis Neurosci* **14:** 1015–1027.

NC Rust, O Schwartz, JA Movshon and EP Simoncelli (2004) Spike-triggered characterization of excitatory and suppressive

stimulus dimensions in monkey V1. *Neurocomputing*, Elsevier.

NC Rust & JA Movshon (2005) In praise of artifice. *Nat Neurosci* **8:** 1647–50.

RR de Ruyter van Stevenينck, GD Lewen, SP Strong & W Bialek (1997) Reproducibility and variability in neural spike trains. *Science* **275:** 1805–1808.

R de Ruyter van Steveninck, WH Zaagman & HAK Mastebroek (1986), Adaptation of transient responses of a movement-sensitive neuron in the visual system of the blowfly Calliphora erythrocephala. *Biol Cybern* **54 :** 223–226.

RR de Ruyter van Steveninck & W Bialek (1988) Real-time performance of a movement sensitive in the blowfly visual system: Information transfer in short spike sequences. *Proc Roy Soc Lond B* **234:** 379–414.

HM Sakai (1992) White-noise analysis in neurophysiology. *Physiol Rev* **72(2):**  491-505.

M Schetzen (1989) The Volterra and Wiener theories of nonlinear systems. Krieger, Malabar.

E Schneidman, MJ Berry 2nd, R Segev & W Bialek (2006) Weak pairwise correlations imply strongly correlated network states in a neural population. *Nature* **440:** 1007-12.

O Schwartz & EP Simoncelli (2001) Natural signal statistics and sensory gain control. *Nat Neurosci* **4:** 819–25.

O Schwartz, EJ Chichilnisky & E Simoncelli (2002) Characterizing neural gain control using spike triggered covariance. *NIPS* **14:** 269–276.

O Schwartz, JW Pillow, NC Rust & EP Simoncelli (2006) Spike-triggered neural characterization. *J Vis* **6:** 484–507.

G Schwartz, S Taylor, C Fisher, R Harris, MJ Berry 2nd (2007) Synchronized firing among retinal ganglion cells signals motion reversal. *Neuron* **55:** 958-69.

TO Sharpee, NC Rust & W Bialek (2004) Analyzing neural responses to natural signals: maximally informative dimensions. *Neural Comput* **16:** 223–50.

TO Sharpee, H Sugihara, AV Kurgansky, SP Rebrik, MP Stryker & KD Miller (2006) Adaptive filtering enhances information transmission in visual cortex. *Nature* **439:** 936–42.

EP Simoncelli & BA Olshausen (2001) Natural image statistics and neural representation. *Annu Rev Neurosci* **24:** 1193–216.

EP Simoncelli, L Paninski, J Pillow & O Schwartz (2004) Characterization of neural responses with stochastic stimuli. In Gazzaniga M (ed), The Cognitive Neurosciences, 3rd ed. MIT Press, Cambridge, MA.

SM Smirnakis, MJ Berry, DK Warland, W Bialek & M Meister (1997) Adaptation of retinal processing to image contrast and spatial scale. *Nature* **386:** 69–73.

G Tkačik & MO Magnasco (2008) Decoding spike timing: the differential reverse-correlation method. *Biosystems* **93:** 90–100.

G Tkačik, JS Prentice, JD Victor & V Balasubramanian (2010) Local statistics in natural scenes predict the saliency of synthetic textures. *Proc Nat'l Acad Sci USA* **107:** 18149–54.

G Tkačik, P Garrigan, C Ratliff, G Milčinski, JM Klein, LH Seyfarth, P Sterling, DH Brainard & V Balasubramanian (2011) Natural images from the birthplace of the human eye. *PLoS ONE* **6:** e20409.

J Touryan, B Lau & Y Dan (2002) Isolation of relevant visual features from random stimuli for cortical complex cells. *J Neurosci* **22:** 10811–8.

W Truccolo, UT Eden, MR Fellows, JP Donoghue & EN Brown (2004) A point process framework for relating neural spiking activity to spiking history, neural ensemble and extrinsic covariate effects. *J Neurophysiol* **93:** 1074–89.

JD Victor & P Johannesma (1986) Maximum-entropy approximations of stochastic nonlinear transfuctions: an extension of the wiener theory. *Biol Cybern* **54:** 289–300.

JD Victor & BW Knight (1979) Nonlinear analysis with an arbitrary stimulus ensemble. *Quart Appl Math* **2:** 113–136.

JD Victor & RM Shapley (1979) The nonlinear pathway of Y ganglion cells in the cat retina. *J Gen Physiol* **74:** 671–689.

JD Victor & RM Shapley (1980) The effect of contrast on the non-linear response of the Y cell. *J Physiol* **302:** 535–547.

N Wiener (1958) Nonlinear problems in random theory. Wiley, New York.

MCK Wu, SV David & JL Gallant (2006) Complete functional characterization of sensory neurons by system identification. *Annu Rev Neurosci* **29:** 477–505.

C Zetzsche & U Nuding (2005) Nonlinear and higher-order approaches to the encoding of natural scenes. *Network* **16:** 191–221.