

# **Tight data-robust bounds to mutual information combining shuffling and model selection techniques**

M.A. Montemurro<sup>1</sup>, R. Senatore<sup>1</sup> and S. Panzeri<sup>1</sup>

<sup>1</sup> *The University of Manchester, Faculty of Life Sciences,  
The Mill, G Floor,  
PO BOX 88, Manchester M60 1QD, UK*

Please address correspondence to S. Panzeri at the above address.

Tel. (44) 161 3063877, Fax (44) 161 3063887, Email: [s.panzeri@manchester.ac.uk](mailto:s.panzeri@manchester.ac.uk)

## Abstract

The estimation of the information carried by spike times is crucial for a quantitative understanding of brain function, but it is difficult because of an upward bias due to limited experimental sampling. We present new progress, based on two basic insights, on reducing the bias problem. First, we show that, by means of a careful application of data shuffling techniques, it is possible to cancel almost entirely the bias of the noise entropy, the most biased part of information. This procedure provides a new information estimator which is much less biased than the standard direct one and has similar variance. Second, we use a non-parametric test to determine whether all the information encoded by the spike train can be decoded assuming a low dimensional response model. If this is the case, the complexity of response space can be fully captured by a small number of easily sampled parameters. Combining these two different procedures we obtain a new class of precise estimators of information quantities, which can provide data-robust upper and lower bounds to the mutual information. These bounds are tight even when the number of trials per stimulus available is one order of magnitude smaller than the number of possible responses. The effectiveness and the usefulness of the methods is tested through application to simulated data and to recordings from somatosensory cortex. This application shows that, even in the presence of strong correlations, our methods constrain precisely the amount of information encoded by real spike trains recorded in vivo.

# 1 Introduction

A recent fundamental insight from system neuroscience is that information about the external sensory world is often encoded by the precise timing of action potential (spikes). The time of individual spikes has been shown to encode precisely and reliably the occurrence of certain stimulus features (Rieke et al., 1996; de Ruyter van Steveninck et al., 1997; Buracas et al., 1998; Panzeri et al., 2001b; DeWeese et al., 2003; Arabzadeh et al., 2005). However, one question which is still unanswered is how individual spikes, either from the same or from a different neuron, combine together to give rise to perception. One fundamental observation is that spike times are correlated: for example, nearby cortical cells tend to fire in synchrony more than expected by chance. The presence of correlations has suggested that they are a fundamental ingredient of the neural code. The computational advantage of such a representation may be that correlations add an additional information channel that can be used to either represent more sensory/behavioral features (Abeles et al., 1993; Dan et al., 1998) or to bind together groups of features (Gray et al., 1989; von der Malsburg, 1999). However, whether or not correlations are a crucial part of the neural code is still highly controversial (Shadlen and Movshon, 1999).

A principled and rigorous way to address how the messages carried by individual spike times are integrated together is to use information theory to quantify and compare different ways and time scales at which spike times may convey information (Rieke et al., 1996; Borst and Theunissen, 1999; Dimitrov and Miller, 2001; Panzeri et al., 2001b). The use of information theory allows an estimate of how reliably stimuli are encoded in single trials, and of which features of the neuronal response, such as independent spikes or the correlations, contribute to stimulus discriminability.

However, a problem with this approach is that quantifying reliably the information conveyed by spike timing often requires the collection of unpractically large samples of data. This is mainly due to correlations: if correlations did not exist, then the statistics of spike times would be completely characterized by the time-dependent firing rate of each neuron. However, one also needs to measure the correlations among all possible groups of spikes. A complete characterization of these correlations requires a number of parameters which are difficult to sample with realistic amounts of neuronal data. Thus, spike timing information measures suffer from a significant sampling bias problem (Panzeri and Treves, 1996).

Generalizing previous work of Reich and colleagues (Reich et al., 2000), we have recently proposed an approach to alleviate the sampling bias problem, by developing data-robust lower bounds to the spike timing information that neglect long-lag stimulus-modulations of correlations (Pola et al., 2005). These bounds can establish if there is information conveyed in spike times above and beyond that conveyed by spike counts. However, these methods cannot be used to test the importance of correlations in coding, because they explicitly neglect a potentially

important part of the correlation structure.

In this paper we overcome this limitation by presenting several methodological advances which lead to a radical improvement of the sampling properties of information measures and provide very tight and data-robust upper- and lower-bounds to spike timing information. These advances also permit to constrain precisely the role of correlations in decoding. Overall, this progress provides the basis for a better determination of the role of correlations in information transmission, and at the same time significantly expands the domain of applicability of information analysis techniques in the analysis of neural signals.

The paper is organized as follows: We first review basic concepts of information theory applications to spike trains; we then discuss how to quantify the importance of correlations in decoding; we then address the sampling properties of these information quantities and provide bounds that are either biased upwards or downwards; we then discuss how to use model selection techniques to give virtually unbiased and tight estimations of information. Lastly we apply the new techniques to real neuronal spike trains recorded from rat somatosensory cortex.

## 2 The information carried by neuronal population responses

We consider a time period of duration  $T$ , associated with a dynamic or static sensory stimulus  $s$  (chosen with probability  $P(s)$  from a stimulus set  $\mathcal{S}$  with  $S$  elements), during which the activity of one neuron is observed<sup>1</sup>. We assume that the spike arrival times are binned with a timing precision  $\Delta t$  and transformed into a sequence of spike counts in each time bin.  $L$  denotes the number of time-bins (i.e.  $T = L\Delta t$ ). The neuronal response is denoted by a one-dimensional array  $\mathbf{r} = \{r(1), r(2), \dots, r(L)\}$ , where  $r(t)$  is the number of spikes emitted by the neuron in the  $t$ -th time-bin. The maximum number of spikes that can be observed in a single time bin in any trial is denoted by  $M$ . (If  $\Delta t$  is very short,  $M$  is 1 and  $r(t)$  is binary). We indicate the response space by  $\mathcal{R}$  ( $\mathcal{R}$  contains  $(M + 1)^L$  elements).

Following Shannon (1948), we write the mutual information  $I(\mathcal{R}; \mathcal{S})$  (often abbreviated as  $I$  in the following) transmitted by the population response about

---

<sup>1</sup>In this paper we consider one neuron only in order to keep notations simpler. However, the generalization to neuronal populations is relatively straightforward and does not present conceptual difficulties. The main step in generalizing to populations is a slight change needed in the definitions of the Markov models described below, see (Pola et al., 2005) for an example of how to carry out this generalization.

the whole set of stimuli as

$$I(\mathcal{R}; \mathcal{S}) = H(\mathcal{R}) - H(\mathcal{R}|\mathcal{S}), \quad (1)$$

where  $H(\mathcal{R})$  and  $H(\mathcal{R}|\mathcal{S})$  are the *response entropy* (stimulus-unconditional) and the *noise entropy* (stimulus-conditional) respectively. They are defined (Cover and Thomas, 1991) as

$$H(\mathcal{R}) = - \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}) \log_2 P(\mathbf{r}), \quad (2)$$

$$H(\mathcal{R}|\mathcal{S}) = - \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}|s) \log_2 P(\mathbf{r}|s). \quad (3)$$

The response entropy quantifies how neuronal responses vary with the stimulus and thus sets the capacity of the spike train to convey information. The noise entropy quantifies the irreproducibility of the neuronal responses at fixed stimulus. Thus, mutual information quantifies how much of the information capacity provided by stimulus-evoked differences in neural activity is robust to the presence of trial-by-trial response variability (de Ruyter van Steveninck et al., 1997). In Eqs. (2) and (3) the summation over  $\mathbf{r}$  is over all possible neuronal responses. The summation over  $s$  is over all possible stimuli.  $P(\mathbf{r}|s)$  is the probability of observing a particular response  $\mathbf{r}$  conditional to stimulus  $s$ . Experimentally,  $P(\mathbf{r}|s)$  is determined by repeating each stimulus on many trials, while recording the neuronal responses. The probability  $P(s)$  is usually chosen by the experimenter.  $P(\mathbf{r}) = \langle P(\mathbf{r}|s) \rangle_s$  is its average across all stimuli (the angular brackets indicate an average over stimuli,  $\langle F(s) \rangle_s \equiv \sum_{s \in \mathcal{S}} P(s) F(s)$ ). We assume that there are enough stimuli in the presented set so that  $P(\mathbf{r})$  (which is computed across all trials to all stimuli) is better sampled than  $P(\mathbf{r}|s)$ . (In practice, this amounts to the requirement that more than a handful of stimuli is presented).

Estimating the information carried by spike times of real neuronal populations is difficult because each stimulus-response probability has to be measured from a limited amount of data. The statistical errors in estimating the response probabilities lead to a downward systematic error (bias) in both noise and response entropy (Miller, 1955).  $H(\mathcal{R})$  depends only on  $P(\mathbf{r})$ , which is sampled across all trials to all stimuli. Under our assumptions, its bias is much smaller than that of  $H(\mathcal{R}|\mathcal{S})$ , which depends on  $P(\mathbf{r}|s)$ . This results in an overall upward bias when estimating mutual information (Panzeri and Treves, 1996). This makes it difficult to estimate the information directly from Eq. (1), especially for long time windows or precise spike time discretizations (large  $L$ ).

### 3 Simplified models of correlation

Having defined the information that neuronal responses transmit about sensory stimuli, we consider how correlations in the responses affect information transmission.

The first step is to define precisely what we mean by correlations. In this paper, when we say that the spike trains are correlated we mean that, for some stimulus  $s$ , the “true” stimulus-response probability  $P(\mathbf{r}|s)$  is different from the probability  $P_{ind}(\mathbf{r}|s)$  obtained if spikes were independent at fixed stimulus. By definition, the independent probability model  $P_{ind}(\mathbf{r}|s)$  is the product of the stimulus-conditional marginal probabilities  $P(r(t)|s)$  of responses in each time-bin  $t$ :

$$P_{ind}(\mathbf{r}|s) = \prod_{t=1}^L P(r(t)|s), \quad (4)$$

Thus, when we refer to correlations we mean correlations at fixed stimulus. These correlations are usually called *noise correlations* (Gawne and Richmond, 1993; Nirenberg and Latham, 2003; Pola et al., 2003). For brevity, in the rest of this paper when we use the term “correlation” we mean “noise-correlation”.

After correlations have been defined, the next step is to characterize how they affect information transmission. Correlations can affect neural information transmission in different ways, both in terms of encoding and downstream decoding of neuronal messages. Here, following previous work (Latham and Nirenberg, 2005; Nirenberg and Latham, 2003), we specifically focus on whether or not correlations must be taken into account to *decode* the neuronal response. We consider a downstream neural system that bases its decoding decisions on the assumption that the spikes are generated by a simplified response model  $P_{simp}(\mathbf{r}|s)$  which neglects certain aspects of the spike train correlation structure (e.g., it considers only correlations between spikes close together in time)<sup>2</sup>. We ask how much information is lost because the decoding operation is performed assuming that responses  $\mathbf{r}$  are generated with  $P_{simp}(\mathbf{r}|s)$  rather than with  $P(\mathbf{r}|s)$ .

The choice of the mathematical form of  $P_{simp}(\mathbf{r}|s)$  will depend on the question that the experimenter wants to address about correlations. If for example one is interested in whether correlations of any form are important for decoding, then one considers how much information is lost when the independent model  $P_{ind}(\mathbf{r}|s)$  is used for decoding. If instead one is interested in the more specific question of whether correlations within a specified time range are important for decoding, then one considers how much information would be lost when using simplified response models  $P_{simp}(\mathbf{r}|s)$  that neglect correlations at time scales outside the specified range. When considering neuronal population recordings, a similar strategy could be used to study the spatial scale at which correlations influence information transmission: in this case  $P_{simp}(\mathbf{r}|s)$  will take into account only correlations among a specific subset of neurons.

---

<sup>2</sup>For example, the downstream system may decode the stimulus using, via Bayes’ rule, a posterior probability based on the simplified model:  $P_{simp}(s|\mathbf{r}) = P(s)P_{simp}(\mathbf{r}|s)/P_{simp}(\mathbf{r})$ .

### 3.1 Assumptions on the simplified models of correlation

Although in the remainder of the paper we will focus on a particular class of simplified response models, in this section it is useful to keep the simplified probability model  $P_{simp}(\mathbf{r}|s)$  as general as possible and spell out what are the minimal requirements to  $P_{simp}(\mathbf{r}|s)$  that are necessary to develop our information theoretic formalism. We will only require that  $P_{simp}(\mathbf{r}|s)$  satisfies two assumptions. These assumptions are listed below, and their importance will be made clear in the rest of the paper.

Before we list the assumptions we note that we will describe our formalism by having in mind simplified models  $P_{simp}(\mathbf{r}|s)$  that depend on much less parameters than  $P(\mathbf{r}|s)$ , and are thus much easier to sample than  $P(\mathbf{r}|s)$ . For example, the simplest possible case of  $P_{simp}(\mathbf{r}|s)$  is a parameter-free probability distribution which is uniformly flat across all times and stimuli. Another example of a model which is simple to sample is  $P_{simp}(\mathbf{r}|s) = P_{ind}(\mathbf{r}|s)$ . In fact, while estimating  $P(\mathbf{r}|s)$  requires an evaluation of  $(M + 1)^L - 1$  parameters for each stimulus  $s$ , estimating  $P_{ind}(\mathbf{r}|s)$  needs only  $ML$  parameters for each stimulus. Although we have in mind very simple models for  $P_{simp}(\mathbf{r}|s)$ , it is important to note that the formalism developed in this paper would be well defined even if  $P_{simp}(\mathbf{r}|s)$  is approximately as complex to sample as  $P(\mathbf{r}|s)$ . In fact, the family of Markov models that we will study in detail below interpolates parametrically from low to high model complexity.

We require that the simplified model  $P_{simp}(\mathbf{r}|s)$  to be used satisfies the following two assumptions:

*Assumption 1.* We require that the method used for transforming  $P(\mathbf{r}|s)$  to  $P_{simp}(\mathbf{r}|s)$  operates separately and independently on the responses conditioned to each stimulus. Thus, we require that the transformation from  $P(\mathbf{r}|s)$  to  $P_{simp}(\mathbf{r}|s)$  is independent of  $P(s)$ , or, of  $P(\mathbf{r}|s')$  and  $P(s')$ , for any  $s' \neq s$ . This property is useful because the resampling (or “shuffling”) techniques that, as detailed in Section 4, can reduce the bias of information estimates, can only be applied to  $P_{simp}(\mathbf{r}|s)$  that, for each stimulus  $s$ , are constructed only from responses collected in response to that stimulus.

*Assumption 2.* We require that, for each stimulus  $s$  with non-zero probability, the simplified response model  $P_{simp}$  satisfies the following condition:

$$\sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}|s) \log_2 P_{simp}(\mathbf{r}|s) = \sum_{\mathbf{r} \in \mathcal{R}} P_{simp}(\mathbf{r}|s) \log_2 P_{simp}(\mathbf{r}|s) \quad (5)$$

Assumption 2 is important to our analysis for three reasons. The first reason is that (taking the point of view that  $0 \log(0)$  is zero and  $c \log(0)$  is ill-defined for any  $c \neq 0$ ) assumption 2 enforces the condition that if, for some  $\mathbf{r}$  and  $s$ ,  $P_{simp}(\mathbf{r}|s)$  is zero, then  $P(\mathbf{r}|s)$  must also be zero. This fact is crucial in the present context because, as we will see in the next Subsection, it ensures that the information theoretic quantities to be introduced below are well defined. The second reason

is that, as also shown in the next Subsection, assumption 2 ensures that we can rewrite the information theoretic quantities in a way that is more easy to sample. The third reason is that assumption 2 is satisfied by all maximum-entropy models constrained to preserve selected features of the full probability model  $P(\mathbf{r}|s)$ . This will be demonstrated in Section 3.3. Since maximum-entropy smoothing is a principled way to fill the unconstrained details of a simplified model, it is useful to ensure that our formalism is applicable to all such models.

### 3.2 Measures of the information lost in decoding with the simplified model

Now we turn to determining how much information is lost when decoding the neural response with a mismatched decoding model  $P_{simp}(\mathbf{r}|s)$  instead of with the true model  $P(\mathbf{r}|s)$ . This problem has been well studied in the information theoretic literature (Merhav et al., 1994; Latham and Nirenberg, 2005). Although this information loss cannot be expressed through a general and simple analytical expression, Latham and Nirenberg (2005) recently derived a simple closed-form expression that is an upper bound to it, as follows:

$$\Delta I_{simp} \equiv D(P(s|\mathbf{r})||P_{simp}(s|\mathbf{r})) \equiv \sum_{\mathbf{r}} P(\mathbf{r}) \sum_s P(s|\mathbf{r}) \log_2 \frac{P(s|\mathbf{r})}{P_{simp}(s|\mathbf{r})} \quad (6)$$

where  $D$  is conditional Kullback-Leibler (KL) distance (see Cover and Thomas (1991), p. 22, Eq. (2.65)). Assumption 2 above ensures that if, for some  $\mathbf{r}$  and  $s$ ,  $P_{simp}(\mathbf{r}|s)$  is zero, then  $P(\mathbf{r}|s)$  must also be zero, and this in turn ensures that  $\Delta I_{simp}$  is a non-negative and non-divergent information-theoretic measure.

An important problem in the practical estimation of  $\Delta I_{simp}$  in the case in which  $P_{simp}(s|\mathbf{r})$  is described by much less parameters than  $P(s|\mathbf{r})$ , is that it is heavily biased, approximately as much as the mutual information  $I$  (Pola et al., 2005). In the next sections, we will show how to reduce the bias problem of  $\Delta I_{simp}$  and thus allow its estimation in practice.

A second quantity of interest is  $I_{LB-simp}$  (see (Pola et al., 2005)), the difference between the mutual information  $I$  and  $\Delta I_{simp}$ :

$$\begin{aligned} I_{LB-simp} &= I - \Delta I_{simp} \\ &= - \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}) \log_2 P_{simp}(\mathbf{r}) + \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}|s) \log_2 P_{simp}(\mathbf{r}|s) \\ &= \chi_{simp}(\mathcal{R}) + \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}|s) \log_2 P_{simp}(\mathbf{r}|s). \end{aligned} \quad (7)$$

where

$$\chi_{simp}(\mathcal{R}) \equiv - \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}) \log_2 P_{simp}(\mathbf{r}) \quad (8)$$



Since  $\Delta I_{simp}$  is non-negative and is an upper bound to the information lost when decoding the neuronal responses with the mismatched response model  $P_{simp}$ ,  $I_{LB-simp}$  has a well defined meaning: it quantifies information that can be decoded by using  $P_{simp}$ . As shown by Pola et al. (2005), this quantity is of practical importance because it is much less biased than the mutual information  $I$ ; therefore it provides a useful data-robust lower bound to the information decodable with  $P_{simp}$ .

A further simplification to both  $\Delta I$  and  $I_{LB}$  can be obtained by making use of assumption 2, which permits to rewrite  $I_{LB-simp}$  and  $\Delta I_{simp}$  as:

$$I_{LB-simp} = \chi_{simp}(\mathcal{R}) - H_{simp}(\mathcal{R}|\mathcal{S}) \quad (9)$$

$$\Delta I_{simp} = H_{simp}(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S}) + H(\mathcal{R}) - \chi_{simp}(\mathcal{R}) \quad (10)$$

where  $H_{simp}(\mathcal{R}|\mathcal{S})$  is the noise entropy of the simple response model:

$$H_{simp}(\mathcal{R}|\mathcal{S}) = - \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r} \in \mathcal{R}} P_{simp}(\mathbf{r}|s) \log_2 P_{simp}(\mathbf{r}|s). \quad (11)$$

The advantage of this rewriting is that now the stimulus-conditional functionals of the simplified model are expressed as the noise entropy of  $P_{simp}(\mathbf{r}|s)$ . This property is important to improving the sampling properties of the information quantities, because, as we will see in the next Section,  $H_{simp}(\mathcal{R}|\mathcal{S})$  can be corrected for limited sampling by very effective techniques (Nemenman et al., 2004), and its presence in the expression for  $\Delta I_{simp}$  will allow us to cancel out the bias of the latter with appropriate procedures.

### 3.3 Maximum entropy models

We now consider in more detail the problem of how to construct simplified correlation models which satisfy the assumptions needed by our information theoretic framework.

In constructing simplified models of correlations, it is natural to ask our model to preserve only some properties of the true probability  $P(\mathbf{r}|s)$ . A way to formalize this is to require our simplified model to satisfy, apart from the usual requirements of non-negativity and normalization to one, a certain number  $m$  of constraints that are also satisfied by  $P(\mathbf{r}|s)$ , as follows:

$$\begin{aligned} P_{simp}(\mathbf{r}|s) &> 0 \\ \sum_{\mathbf{r}} P_{simp}(\mathbf{r}|s) &= 1 \\ \sum_{\mathbf{r}} P_{simp}(\mathbf{r}|s) g_i(\mathbf{r}) &= \sum_{\mathbf{r}} P(\mathbf{r}|s) g_i(\mathbf{r}) \quad i = 1, \dots, m \end{aligned} \quad (12)$$

where  $g_i(\mathbf{r})$  are arbitrary functions on  $\mathcal{R}$ <sup>3</sup>. Once the constraints in Eq. (12) have been chosen, it is then desirable to simplify the response model by removing

---

<sup>3</sup>The functions  $g_i(\mathbf{r})$  could in principle be different for each stimulus conditional distribution.

all types of correlation in the data apart from those enforced by the features preserved from the original distribution.

A principled way to choose a  $P_{simp}(\mathbf{r}|s)$  that satisfy the constraints in Eq. (12) and adds no further relationship between the data is to choose  $P_{simp}(\mathbf{r}|s)$  as the distribution with the maximum entropy allowed by our constraints. This maximum entropy distribution is unique and has the following expression (Cover and Thomas, 1991):

$$P_{simp}(\mathbf{r}|s) = \exp \left\{ \lambda_0 - 1 + \sum_{i=1}^m \lambda_i g_i(\mathbf{r}) \right\} \quad (13)$$

where the parameters  $\lambda_0, \lambda_i, \dots, \lambda_m$  are fixed (independently for each conditional distribution to each stimulus) so as to satisfy the constraints in Eq. (12). The maximum entropy distribution is in some way the most reasonable choice of simplified model of correlations given the constraints: to choose a distribution with lower entropy would correspond to assume some additional structure that we do not know; to choose one with a higher entropy would necessarily violate the constraints that we wish to enforce.

It is important to note that any maximum-entropy simplified model of the form in Eq. (13) that satisfies constraints of Eq. (12) is a suitable simplified model for our analysis; in fact any such maximum-entropy model satisfies by construction the two assumptions of our formalism. In particular, by using Eq. (13), Eq. (5) of assumption 2 becomes:

$$\sum_{\mathbf{r}} P_{simp}(\mathbf{r}|s) \left\{ \lambda_0 - 1 + \sum_{i=1}^m \lambda_i g_i(\mathbf{r}) \right\} = \sum_{\mathbf{r}} P(\mathbf{r}|s) \left\{ \lambda_0 - 1 + \sum_{i=1}^m \lambda_i g_i(\mathbf{r}) \right\} \quad (14)$$

which is obviously satisfied if  $P_{ind}(\mathbf{r}|s)$  meets the constraints set by Eq. (12). Another demonstration of the relationship between assumption 2 and the maximum-entropy principled is reported in Appendix B.

An important class of maximum entropy distributions is made of the distributions that preserve some marginal probabilities of  $P(\mathbf{r}|s)$ , such as e.g. the independent model  $P_{ind}(\mathbf{r}|s)$  of Eq.(4), the Markov models considered in the next subsection, and the hierarchical probability models of Amari (2001). Models preserving marginals are obtained from Eq. (13) by constraining  $P_{simp}$  and  $P$  to have equal sum on a number of subsets  $\mathcal{A}_i, \mathcal{B}_i, \dots$  of the responses space  $\mathcal{R}$  (each subset corresponding to the responses that have to be summed to compute the marginal probability to be preserved). This corresponds to choose, for each subset, a function  $g_i(\mathbf{r})$  in Eqs. (12,13) with value one on the subset under consideration and zero elsewhere.

For example, the independent model of Eq. (4) is obtained from Eq. (13) by partitioning the response space into the disjoint union of the “marginal” subsets  $\mathcal{A}_1, \mathcal{A}_2$  (the subset of all responses with respectively 1 or zero spikes in time bin 1),  $\mathcal{B}_1, \mathcal{B}_2$  (the subset of all responses with respectively 1 or zero spikes in time

bin 2), and so on for all time bins, and then by enforcing  $P_{simp}$  and  $P$  to have equal sum on each of the subsets.

More complex maximum entropy models can be obtained by extending the procedure to constrain  $P_{simp}$  also on intersections of subsets. For example, Markov models of order 1 can be obtained by constraining the sum of  $P_{simp}$  not only on the intersections of the above marginal subsets corresponding to each time bin, but also on the pairwise intersections corresponding to adjacent time bins. Amari’s hierarchical models of purely pairwise interactions (Amari, 2001) can be obtained by the more strict requirement of constraining the sum of  $P_{simp}$  on all the pairwise intersections of marginal subsets, not only to that corresponding to adjacent time bins.

All of these transformations can be applied recursively. The recursion leads to constructing Markov chains of arbitrary length, as well as hierarchical models containing higher order interactions (Amari, 2001). It is interesting to note that in this way one can prove that all suffix trees models satisfy our assumptions and therefore are valid choices of  $P_{simp}$ . In fact, suffix tree models can be constructed with this recursive partitioning, by constraining  $P_{simp}$  on a smaller number of intersections than the Markov model of corresponding length. This observation helps to understand the relationship between the work presented here and the recent work of London et al. (2002) and Kennel et al. (2005), which use suffix tree models to estimate entropy rates.

In summary, in this section we have established that all maximum-entropy models satisfy our assumptions and can thus be used to estimate information with our formalism. These models include Markov chains, hierarchical distributions and suffix tree models.

We have also provided a general and explicit way to construct such simplified models from the data through Eqs. (12,13). This construction can be successfully applied whatever the statistics of neuronal firing described by the response probability  $P(\mathbf{r}|s)$ . In fact, for any given choice of the constraints in Eq. (12), the maximum entropy model fulfilling these constraints will automatically satisfy (by construction) our assumptions, whatever the form of  $P(\mathbf{r}|s)$ . An important implication of this result is that, in practice, our assumptions *will not restrict* in any way the applicability of the method only to datasets with specific statistical properties.

### 3.4 Markov models

Despite the generality of the above constructions of  $P_{simp}$ , we will illustrate and develop the main idea behind our formalism by focusing on a specific class of simplified correlation models: the Markov models with finite-memory. We choose to illustrate our ideas using this particular class of maximum entropy simplified model because (i) by tuning the order of the Markov process, we can vary parametrically the complexity of the model and thus illustrate clearly how the

sampling behavior of the information theoretic functional depends on the number of parameters describing the simplified model, and (ii) these models are easy to construct and apply to data.

A neurophysiological motivation of the use of Markov models stems from the fact that in many neural systems correlations are significant only between spikes that are separated by a short time lag, in the range 1-15 ms (Gray et al., 1989; Brosch et al., 1997; Dan et al., 1998; Nirenberg et al., 2001; Golledge et al., 2003). In such cases, to preserve the whole information, it is sufficient to take into account *only* correlations extending over a short lag. Thus, one can approximate the real probability of current response  $r(t)$  given the past firing with a finite-memory Markov model that looks back only  $q$  time-steps, as follows:

$$\begin{aligned} P_q(\mathbf{r}|s) &= P(r(1)|s) \prod_{t=2}^L P(r(t)|r(t-q), \dots, r(t-1); s), \quad \text{if } q = 1, \dots, L-1, \\ P_0(\mathbf{r}|s) &= P_{ind}(\mathbf{r}|s), \quad \text{if } q = 0. \end{aligned} \quad (15)$$

The probability conditional on the response in the previous time steps at fixed stimulus in the above equation can be computed from the experimental probabilities via:

$$P(r(t)|r(t-q), \dots, r(t-1); s) = \frac{P(r(t-q), \dots, r(t-1), r(t)|s)}{P(r(t-q), \dots, r(t-1)|s)}. \quad (16)$$

where  $P(r(t-q), \dots, r(t-1), r(t)|s)$  and  $P(r(t-q), \dots, r(t-1)|s)$  are marginal distributions of the full model  $P(\mathbf{r}|s)$ , computed by integrating away the dependence on all the response variables which do not enter in their argument. This simple procedure to construct the Markov model is equivalent to the maximum-entropy one described above.

Markov models interpolate parametrically between the independent model ( $q = 0$ ) and the full probability model (obtained for  $q = L-1$ , because  $P_{L-1}(\mathbf{r}|s) = P(\mathbf{r}|s)$ ).  $P_q(\mathbf{r}|s)$  preserves all correlations extending up to  $q$  time-bins in the past, and it neglects all correlations of range longer than  $q$ . Thus, it is a perfect description of neuronal firing if correlations extend to a lag shorter than or equal to  $q$  time bins.

The information-theoretic probability functionals corresponding to the choice  $P_{simp}(\mathbf{r}|s) = P_q(\mathbf{r}|s)$  will be indicated by a subscript  $q$  in place of the subscript *simp*. For completeness, their expression is reported below:

$$I_{LB-q} = \chi_q(\mathcal{R}) - H_q(\mathcal{R}|\mathcal{S}) \quad (17)$$

$$\Delta I_q = H_q(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S}) + H(\mathcal{R}) - \chi_q(\mathcal{R}) \quad (18)$$

where  $H_q(\mathcal{R}|\mathcal{S})$  is the noise entropy of the simple response model:

$$H_q(\mathcal{R}|\mathcal{S}) = - \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r} \in \mathcal{R}} P_q(\mathbf{r}|s) \log_2 P_q(\mathbf{r}|s) \quad (19)$$

and  $\chi_q(\mathcal{R})$  is

$$\chi_q(\mathcal{R}) = - \sum_{\mathbf{r} \in \mathcal{R}} P(\mathbf{r}) \log_2 P_q(\mathbf{r}) \quad (20)$$

## 4 Bias cancelations obtained by shuffling the responses

As discussed above, mutual information has been broken down into two terms,  $I_{LB-simp}$  and  $\Delta I_{simp}$ , with radically different sampling properties, the former easy to sample and the latter very difficult to sample. Here we examine in detail the sampling behavior of  $\Delta I_{simp}$  and show how to reduce its bias dramatically without increasing its variance. Since  $\Delta I_{simp}$  is the most biased part of  $I$ , this in turn will also improve the sampling properties of the mutual information. Since  $I$ ,  $I_{LB-simp}$  and  $\Delta I_{simp}$  consist of four quantities  $H(\mathcal{R}|\mathcal{S})$ ,  $H_{simp}(\mathcal{R}|\mathcal{S})$ ,  $H(\mathcal{R})$  and  $\chi_{simp}(\mathcal{R})$ , the relative sampling properties of  $I$ ,  $I_{LB-simp}$  and  $\Delta I_{simp}$  can be established by considering the sampling properties of the above four quantities.

### 4.1 The independent decoder - ignoring all correlations

For clarity, we start by considering in detail the sampling properties of the quantities corresponding to the  $q = 0$ , independent probability model  $P_{simp}(\mathbf{r}|s) = P_{ind}(\mathbf{r}|s) = P_0(\mathbf{r}|s)$ . The corresponding information theoretic functionals are  $I_{LB-0}$  and  $\Delta I_0$ , taken from Eqs. (19) and (20) with  $q = 0$ .

#### 4.1.1 Uncorrected estimators

There are two relevant aspects of the sampling properties of a probability functional: its bias and its variance. We will consider the bias first, and the variance later. The bias of a functional  $F(P)$  of the probabilities  $P$  is defined as the difference between  $\langle F(P_N) \rangle_N$ , the ensemble-averaged value of the functional computed from the probability distributions  $P_N$  empirically obtained from  $N$  trials, and the true value of the functional  $F(P)$  computed with the true probability distribution  $P$ . Thus, the bias is a systematic error that cannot be eliminated just by averaging.

We illustrate the sampling behavior of the entropy functionals by computing them on a set of realistically simulated neuronal data, generated as follows. We simulated the spiking response of one neuron in somatosensory cortex to 49 different stimuli consisting of sinusoidal whisker vibrations with different amplitude and frequency (Arabzadeh et al., 2003; Arabzadeh et al., 2004). We simulated neuronal responses over a 0-50 ms post-stimulus time window. We then digitized, with time precision  $\Delta t$  equal to 5ms, these responses into  $L = 10$  binary words. The simulations were performed using a Markov process with order  $q = 3$ ,

with all the marginal probabilities of order 3 or less and the transition probabilities taken from real responses of a cortical somatosensory neuron<sup>4</sup>. The spike train simulated in this way retains a faithful description of all the real neuronal marginal distributions and of their correlations up to  $q = 3$  time bins.

In Fig. 1A we report the sampling behavior of the four quantities  $\chi_0(\mathcal{R})$ ,  $H(\mathcal{R})$ ,  $H_0(\mathcal{R}|\mathcal{S})$ , and  $H(\mathcal{R}|\mathcal{S})$  as a function of the number of trials per stimulus. For each simulation with a fixed number of trials, we computed the “plug-in” value of the functional by just plugging into their equations the empirical estimates of the probability (without application of any bias correction procedure) and then averaging all simulations with the same number of trials.

We first consider the noise entropy  $H(\mathcal{R}|\mathcal{S})$ . Fig. 1A shows that this is by far the most downward-biased of the four functionals considered. This is because it requires the simultaneous measure of all  $P(\mathbf{r}|s)$  to all stimuli. To understand better its sampling behavior, it is useful to find analytical approximations to the bias. These can be easily derived in the *asymptotic sampling regime*. The latter is defined as the case in which the number of trials  $N_s$  to each stimulus  $s$  is so large each response bin with non-zero probability is observed many times,  $N_s P(\mathbf{r}|s) \gg 1$  (Panzeri and Treves, 1996). In this asymptotic sampling regime, the bias of  $H(\mathcal{R}|\mathcal{S})$  (and similarly, of all other functionals) can be expanded in inverse powers of  $1/N$  ( $N$  being the total number of trials across all stimuli), as follows:

$$\text{Bias}[H(\mathcal{R}|\mathcal{S})] \approx \frac{C_1}{N} + \frac{C_2}{N^2} + \dots, \quad (21)$$

The leading term in  $1/N$  has a simple analytical expression:

$$\text{Bias}[H(\mathcal{R}|\mathcal{S})] \approx -\frac{1}{2N \ln 2} \sum_s (\tilde{R}(s) - 1), \quad (22)$$

where  $\tilde{R}(s)$  denotes the number of “relevant” responses of the stimulus conditional response probability distribution  $P(\mathbf{r}|s)$ , i.e. the number of different responses  $\mathbf{r}$  with non-zero probability of being observed when stimulus  $s$  is presented (Panzeri and Treves, 1996).  $\tilde{R}(s)$  is of order  $(M + 1)^L$  for each stimulus. Thus, it follows that for the bias in Eq. (22) to be small,  $N$  should be much bigger than  $S \times (M + 1)^L$ .

Let us now consider  $H_0(\mathcal{R}|\mathcal{S})$ . It is apparent from Fig. 1A that it is almost unbiased. The reason is that it can be expressed as the sum of single-bin entropies

---

<sup>4</sup>For each stimulus condition  $s$ ,  $N_s$  binary spike trains were generated with a  $q$ -order Markov model as follows (see Eq. (16)): The response in the first bin was assigned to be a 1 or a 0 (spike or no spike) according to  $P(r(1)|s)$ , the latter being computed from the real data. Responses in successive time bins were generated one after the other one by one using the corresponding transition probabilities. For instance, the response at bin  $k$  was generated according to the real-data probability  $P(r(k)|r(k-j), \dots, r(k-1); s)$ , where  $j \leq \min(q, k-1)$ .

(see Pola et al. (2005)). As a consequence, its bias is  $\approx ML/2N \log(2)$ , and is thus much smaller than that of  $H(\mathcal{R}|\mathcal{S})$ .

Fig. 1A shows that the noise entropy  $H(\mathcal{R})$  is considerably simpler to sample than  $H(\mathcal{R}|\mathcal{S})$ . The reason is that, since  $H(\mathcal{R})$  depends only on  $P(\mathbf{r})$ , its bias is approximately  $S$  times smaller than the bias of  $H(\mathcal{R}|\mathcal{S})$ . This is an advantage when many different stimuli are presented.

Fig. 1A shows that the bias of  $\chi_0(\mathcal{R})$  is much smaller than the bias of  $H(\mathcal{R})$ . The reason is as follows. Bias arises from the logarithmic form of entropy functionals. The log in  $\chi_0(\mathcal{R})$  depends on  $P_0(\mathbf{r})$ . Since  $P_0(\mathbf{r})$  is better sampled than  $P(\mathbf{r})$ ,  $\chi_0(\mathcal{R})$  has less bias than  $H(\mathcal{R})$ , whose log depends on  $P(\mathbf{r})$ .

Now we study how the properties of the four functionals combine together to give rise to the sampling properties of  $I(\mathcal{R};\mathcal{S})$ ,  $I_{LB-0}$  and  $\Delta I_0$ .

The bias of the mutual information  $I(\mathcal{R};\mathcal{S})$  is the difference between the biases of  $H(\mathcal{R})$  and  $H(\mathcal{R}|\mathcal{S})$ . As the most biased term is  $H(\mathcal{R}|\mathcal{S})$ , the mutual information is upward biased (Panzeri and Treves, 1996).

The bias of  $I_{LB-0}$  (Eq. (17) with  $q = 0$ ) is the difference between the biases of  $\chi_0(\mathcal{R})$  and  $H_0(\mathcal{R};\mathcal{S})$ . As both these quantities are virtually unbiased, so is  $I_{LB-0}$ .

Let us now consider the bias of  $\Delta I_0$  (Eq. (18) with  $q = 0$ ). When many stimuli are presented, the contribution of the  $\chi_0(\mathcal{R})$  and  $H(\mathcal{R})$  is only very mildly biased and determined by the bias of  $H(\mathcal{R})$ . Thus, most of the bias comes from the stimulus-conditional term  $H_0(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S})$ . Since  $H(\mathcal{R}|\mathcal{S})$  is very strongly biased downward and  $H_0(\mathcal{R}|\mathcal{S})$  is essentially unbiased, the two biases do not cancel out and as result  $\Delta I_0$  is biased upward and behaves like  $-H(\mathcal{R}|\mathcal{S})$ .

An important practical problem is how to reduce the bias of  $H_0(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S})$  and thus improve the sampling properties of  $\Delta I_0$ . A solution to this problem is to compute  $H_0(\mathcal{R}|\mathcal{S})$  not only directly from the single bin marginal probability as in Eq. (19), but by randomly *shuffling* the responses and then computing their entropy. In the  $q = 0$  case considered here, we can generate a new set of shuffled responses to stimulus  $s$  by randomly permuting, for each time bin, the order of trials collected in response to the stimulus  $s$  considered, and then joining together the shuffled responses in different time bins into a response vector  $\mathbf{r}_{0-sh}$ . This shuffling operation leaves each single-time-bin marginal probability unchanged (because responses in each bin are just randomly permuted), while destroying any within-trial correlation between different time bins. We define  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  as the noise entropy of the shuffled distribution:

$$H_{0-sh}(\mathcal{R}|\mathcal{S}) = - \sum_{s \in \mathcal{S}} P(s) \sum_{\mathbf{r}_{sh} \in \mathcal{R}} P_{0-sh}(\mathbf{r}|s) \log_2 P_{0-sh}(\mathbf{r}|s). \quad (23)$$

where  $P_{0-sh}(\mathbf{r}|s)$  is the distribution of response values obtained from  $\mathbf{r}_{0-sh}$ . The asymptotic large  $N$  value of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  is the same of that of  $H_0(\mathcal{R}|\mathcal{S})$ , but its scaling with the number of trials is much different. This is shown in Fig. 1A. Unlike  $H_0(\mathcal{R}|\mathcal{S})$ ,  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  scales with  $N$  approximately as  $H(\mathcal{R}|\mathcal{S})$ , but with

a slightly more negative slope (i.e. slightly more downward bias) as the number of trials decreases. The fact that the bias of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  are of similar size can be intuitively understood from the fact that  $P_{0-sh}(\mathbf{r}|s)$  is sampled with the same number of trials as  $P(\mathbf{r}|s)$  from responses with the same dimensionality. To understand why  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  is *more biased downwards* than  $H(\mathcal{R}|\mathcal{S})$ , we computed the bias of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  in the “asymptotic sampling” regime. We found that the asymptotic bias of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  has the same expression as that of  $H(\mathcal{R}|\mathcal{S})$  in Eq. (22), after replacing  $\tilde{R}(s)$  with  $\tilde{R}_{0-sh}(s)$ , the number of bins relevant to  $P_{0-sh}(\mathbf{r}|s)$ . Since  $P_0(\mathbf{r}|s) = 0$  implies  $P(\mathbf{r}|s) = 0$ , and since the shuffled responses are generated according to  $P_0(\mathbf{r}|s)$ , then it must be that  $\tilde{R}_{0-sh}(s) \geq \tilde{R}(s)$ . Thus, at the leading order in the asymptotic regime,  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  is never less downward biased than  $H(\mathcal{R}|\mathcal{S})$ <sup>5</sup>. Therefore, if we are in the asymptotic sampling regime (or at least the number of trials is enough to make the asymptotic equations a decent estimate of the size of the bias) and if  $\tilde{R}_{0-sh}(s)$  and  $\tilde{R}(s)$  are roughly similar,  $H_{0-sh}(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S})$  will either (i) have a leading-order  $1/N$  bias term which is negative and arises from a partial cancelation of roughly similar leading bias terms of two entropies, or (ii) (in case  $\tilde{R}_{0-sh}(s) = \tilde{R}(s)$ ) it will have a very small bias, only at the order  $1/N^2$  or higher.

The sampling behavior of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  suggests a simple strategy to reduce the bias of  $\Delta I_0$ : use  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  instead of  $H_0(\mathcal{R}|\mathcal{S})$ . We call this estimator  $\Delta I_{0-sh}$ :

$$\Delta I_{0-sh} = H_{0-sh}(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S}) + H(\mathcal{R}) - \chi_0(\mathcal{R}) \quad (24)$$

Simulations in Fig. 1B show that  $\Delta I_{0-sh}$  is much less biased than  $\Delta I_0$ , but it converges to the same asymptotic large- $N$  value. The biases of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  almost cancel each other, leaving an overall small negative bias<sup>6</sup>.

The good sampling properties of  $\Delta I_{0-sh}$  suggest a new, alternative, way to estimate mutual information, as follows:

$$I_{sh} = I_{LB-0} + \Delta I_{0-sh} \quad (25)$$

Since  $I_{LB-0}$  is virtually unbiased and  $\Delta I_{0-sh}$  is mildly biased downward,  $I_{sh}$  is also mildly biased downward. This is confirmed by the simulation results in

---

<sup>5</sup>There are different types of resampling methods that have been used to generate surrogate datasets for the validation of information theoretic calculations (Johnson et al., 2001). It is important to note that the fact that  $\tilde{R}_{0-sh}(s) \geq \tilde{R}(s)$  holds for the resampling procedure “without replacement” used here to generate the surrogate data, but it does not necessarily hold for other resampling techniques “with replacement”.

<sup>6</sup>The idea of using shuffling to cancel out the biases in the stimulus-dependent part of  $\Delta I$  has been first proposed by Nirenberg et al. (2001). However the fact that  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  is more downward biased than  $H(\mathcal{R}|\mathcal{S})$ , and thus  $\Delta I_0$  computed in this way is mildly downward biased, has not been reported before.



Fig. 1B: using  $I_{sh}$  compares very favorably with computing mutual information directly as  $I$  in Eq. (1).

In summary, we now have two classes of estimators of  $\Delta I_0$  and therefore of  $I$ : the direct estimator  $\Delta I_0$  and  $I$  in Eqs. (6) and (1), which are strongly biased upward, and the shuffled estimators  $\Delta I_{0-sh}$  and  $I_{sh}$ , which are mildly biased downward. The difference in the sign of the bias of the two different estimators is very important in practice. In fact, computing  $\Delta I_0$  and  $I$  with both upward- and downward-biased algorithms provides a mean to bound both from above and from below. This can greatly improve the confidence in estimates obtained from experimental data. In the rest of this paper we will devote our attention on how to make both upper and lower bounds more tight than they are in Fig. 1.

#### 4.1.2 Effect of different bias correction methods

The above shows the behavior of the probability functionals when estimated by plugging-in the probabilities directly estimated from the data. However, these estimates can be improved by using available bias correction techniques. In this section we apply different bias correction techniques to the probability functionals in order to understand how best to evaluate each functional from a limited number of data.

We first evaluated the performance of a quadratic extrapolation procedure (Strong et al., 1998) performed computing the quantities from fractions of the data available, and then fitting the resulting data-scaling behavior to a quadratic function of  $1/N$  as in Eq. (21). This procedure should work well in the asymptotic regime where the number of trials  $N$  is large<sup>7</sup>. Simulations suggest that, in practice, a good bias correction from quadratic extrapolations or other asymptotic requires the number of trials per stimulus to be at least as big as the number of possible responses  $R$  (Panzeri and Treves, 1996).

Results are shown in Fig. 2A-B. Fig. 2A shows that  $\chi_0(\mathcal{R})$  after the quadratic extrapolation procedures become completely unbiased, even with as little as 32 trials per stimulus. After applying the quadratic extrapolation,  $H_0(\mathcal{R}|\mathcal{S})$  also becomes almost unbiased even at 32 trials per stimulus. This can be understood by remembering that  $H_0(\mathcal{R}|\mathcal{S})$  can be expressed as a sum of single-bin entropies, and the sampling of each single bin entropy can be considered to be in asymptotic regime even for small number of trials. Thus, the quadratic extrapolation performs very well on  $H_0(\mathcal{R}|\mathcal{S})$ . As a result, the extrapolated  $I_{LB-0}$  estimator

---

<sup>7</sup>We also considered the performance of other bias correction methods developed to work in the asymptotic regime, such as computing analytically the coefficients of the  $1/N$  expansion of the bias as a function of probabilities, (Panzeri and Treves, 1996; Pola et al., 2005). As we found no overall increase in performance with respect to the quadratic extrapolation, we decided to work with the latter, which is easier to implement.

(Fig. 2B) becomes almost unbiased even with as little as 32 trials per stimulus.

Fig. 2A shows that the response entropy  $H(\mathcal{R})$  is still mildly biased after correction, requiring at least 128 trials per stimulus for good estimation. In contrast, the noise entropies  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  remain highly biased even after the extrapolation procedure, requiring at least  $2^{10} = 1024$  trials per stimulus for unbiased estimation. The reason is that  $P(\mathbf{r}|s)$  and  $P_{0-sh}(\mathbf{r}|s)$  are high dimensional, and thus the number of trials needed to get them into the asymptotic sampling regime is much higher. As a consequence,  $\Delta I_0$  is still substantially biased upwards and applying a quadratic extrapolation procedure is still not enough: we still need at least  $2^{10}$  trials per stimulus for unbiased estimation (Fig. 2B). However, the good news from Fig. 2B is that the behavior of  $\Delta I_{0-sh}$  is better than the one of  $\Delta I_0$ . Because of the partial bias cancelation in  $H_{0-sh}(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S})$ , the quadratic extrapolation is able to remove most of the downward bias of  $\Delta I_{0-sh}$  from at least 128 trials/stimulus. Similarly, estimating mutual information as  $I_{sh}$  (rather than directly as  $I$  in the mutual information definition in Eq. (1)), leads to a much less biased estimation (Fig. 2B).

In summary, we simulated data with  $L = 10$  time bins and  $2^{10}$  possible responses, and we use the quadratic extrapolation bias correction. Estimating  $\Delta I_0$  and  $I$  directly requires approx.  $2^{10}$  trials per stimulus, of the order of the size of the response space. However, estimating the very same quantities through the  $\Delta I_{0-sh}$  and  $I_{sh}$  shuffling procedure required only  $2^7$  trials per stimulus, 8 times smaller than the size of the responses space. To check how these results scaled with  $L$ , we simulated data with the same procedure as in Fig. 1 and 2, but using different sizes of post-stimulus windows (ranging from  $L = 8$  to  $L = 14$ ). The results were always comparable to the  $L = 10$  case: Estimating  $\Delta I_0$  and  $I$  directly requires approx. a number of trials per stimulus of the order of the size of the response space. Estimating  $\Delta I_{0-sh}$  and  $I_{sh}$  shuffling procedure required only a number of trials per stimulus 8 times smaller than the size of the responses space (results not shown).

We now consider the effect of using a more sophisticated Bayesian entropy bias correction proposed by Nemenman et al. (2004). This procedure (called NSB) was designed to operate beyond the asymptotic regime<sup>8</sup>, and is briefly reviewed in Appendix A. As shown in (Nemenman et al., 2004), it works well even when data are scarce and the response space is high-dimensional, so that each response only is observed in no more than handful of trials. Fig. 2C and 2D report the values of the functional corrected with the NSB procedure and show that the NSB method in general performs well. The response entropy  $H(\mathcal{R})$

---

<sup>8</sup>We also considered the performance of the procedure of Paninski (2003), which also aims at working beyond the asymptotic sampling regime. We found that, on the present simulated data and information theoretic quantities, the Paninski (2003) procedure helped to reduce the bias but performed not as well as the NSB and the quadratic extrapolation. Thus we omitted the presentation of its results.

(Fig. 2C) is even better behaved than with the quadratic correction and is now essentially unbiased, requiring only 32 trials per stimulus for good estimation. The noise entropies  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  are also better evaluated with the NSB procedure than with the quadratic extrapolation procedure. However, both  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  remain still substantially biased. As a consequence, Fig. 2D shows that  $\Delta I_0$  and  $I$  are still substantially biased upwards, requiring at least  $2^9$  trials per stimulus for unbiased computation. However, the unbiased computation of  $\Delta I_{0-sh}$  and of  $I_{sh}$  requires only  $2^6$  trials per stimulus, because the residual biases of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $-H(\mathcal{R}|\mathcal{S})$  cancel out almost exactly. Estimating mutual information as  $I_{sh}$  (Eq. (25)) works much better than estimating  $I$  directly. However, it should be noted that, as a result of applying the NSB procedure, the residual bias for  $I_{sh}$  is not always negative. Thus, using the Nemenman et al. (2004) procedure to estimate  $I_{sh}$  gives an almost unbiased estimator also for small numbers of trials, but not necessarily a lower bound to the true information (as it instead happens when correcting  $I_{sh}$  with the quadratic extrapolation).

The only entropy term that performed worse when corrected with the NSB procedure is the independent entropy  $H_0(\mathcal{R};\mathcal{S})$ . As discussed in Appendix A, the reason is that the Nemenman et al. (2004) procedure is tailored to work for high-dimensional response spaces, and thus its assumptions do not work well for  $H_0(\mathcal{R};\mathcal{S})$ , which is essentially a sum of single-bin entropies. This problem is particularly serious for low firing rates (see Appendix A).

Lastly, we note that the procedure of Nemenman et al. (2004) has been developed so far only for entropy quantities. Therefore it cannot be applied to  $\chi_0(\mathcal{R})$ . Thus,  $\chi_0(\mathcal{R})$  will be always corrected with the extrapolation procedure which, as demonstrated before, works extremely well for this quantity.

It is now useful to come back to discuss why assumption 2 on the probability model, the one that enabled us to express  $I_{LB-simp}$  and  $\Delta I_{simp}$  as in Eq. (9) and (10), is very important to improve the sampling properties. There are two reasons behind it. The first is that the term outside the logarithm in the stimulus-conditional part of  $I_{LB-simp}$  now is  $P_{simp}(\mathbf{r}|s)$  rather than  $P(\mathbf{r}|s)$ , and this is expected to reduce statistical fluctuations. The second, and more important reason, is that now the stimulus-dependent part of  $\Delta I_{LB-simp}$  can be expressed as a difference of entropies. This has a big impact on the ways the sampling bias of this quantity can be corrected. In fact, the sampling bias correction techniques of Nemenman et al. (2004) and the shuffling bias relationships used above both depend crucially on being able to express the stimulus-conditional part of  $I_{LB-simp}$  and of  $\Delta I_{simp}$  entirely in terms of entropies.

#### 4.1.3 Variance of shuffling-based estimators

We now consider whether the reduction in bias of the shuffled estimators  $\Delta I_{0-sh}$  and  $I_{sh}$  comes at the expense of an increase in their variance. Eq. (25) indicates that the variance of  $I_{sh}$  is largely determined by that of  $\Delta I_{0-sh}$ , its worst sam-

pled component. The most biased terms in the computation of  $\Delta I_{0-sh}$  are the stimulus-conditional entropies  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$ , with their bias almost canceling out (see Eq. (26)). The amount of variance of  $\Delta I_{0-sh}$  will depend on the degree of correlation (across independent realizations with the same number of trials) between the values of  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$ . If these two quantities were positively correlated, then their statistical fluctuations would tend to cancel out, and the resulting variance of  $\Delta I_{0-sh}$  would be under control. If instead  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  were either uncorrelated or negatively correlated across different realizations, then their statistical fluctuations would either not cancel out or even increase, thus making the variance of  $\Delta I_{0-sh}$  larger. The scatter plot in Fig. 3A shows that  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  are strongly positively correlated when taken from the same realization of simulations. (The data were simulated exactly as in Fig. 2 and were obtained using 128 trials per stimulus in each simulation). The fact that  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$  are correlated stems from the fact that fluctuations in the values of single bin marginal probabilities have a major impact on the fluctuations of entropy values (Schultz and Panzeri, 2001). These fluctuations are reflected with the same sign in both  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$ . Thus, the correlation between the entropies ensures that the variances of  $I_{sh}$  and  $I(\mathcal{R};\mathcal{S})$  remain comparable. This is clearly demonstrated in Fig. 3B (quadratic extrapolation correction procedure used).

To better understand the importance of the correlation of  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  across different realizations, we paired values of  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  taken from randomly chosen realizations of the simulation. The scatter plot of these randomly paired entropies is reported in Fig. 3C. Computing information  $I_{sh}$  from the randomly paired entropies  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  leads to a very considerable increase in invariance of the information  $I_{sh}$  (results reported in Fig. 3D).

Thus, we conclude that the reduction of bias in estimating information with  $I_{sh}$  rather than directly with  $I$  does not come at the expense of increased variance, because the computation of  $I_{sh}$  involves a cancelation of fluctuations between similarly biased terms.  $I_{sh}$  is a very efficient estimator of information both in terms of bias and of variance.

#### 4.1.4 Performance of shuffling-based estimators in the presence of strong correlation

An important question is how the bias properties of  $H(\mathcal{R}|\mathcal{S})$ ,  $H_{0-sh}(\mathcal{R}|\mathcal{S})$ , and  $I_{sh}$  are affected by the strength of correlations between spikes. As it was discussed above, when the responses are shuffled the new value of relevant responses,  $\tilde{R}_{0-sh}(s)$  will be equal or larger than the original  $\tilde{R}(s)$ . Thus, because of Eq. (22), if correlations are not strong enough to induce radical differences of shape and support between  $P(\mathbf{r}|s)$  and  $P_{ind}(\mathbf{r}|s)$ ,  $\Delta I_{0-sh}$  is expected to have a downward bias which is much smaller in absolute magnitude than the upward bias of  $\Delta I_0$ .

The previously shown simulations (based on a correlation structure taken from a real cortical neuron) confirmed this expectation and suggest that  $\Delta I_{0-sh}$  and  $I_{sh}$  are only very mildly biased in realistic neurophysiological conditions.

However, this argument also suggests that the magnitude of the downward bias may be affected by the overall strength of correlation. High correlation tends to favor certain patterns in the response and thus decrease the number of possible responses. The stronger the correlations, the larger the number of new response words that may appear after shuffling and the larger the resulting downward bias. If correlations are strong enough to make the number of shuffled relevant responses  $\tilde{R}_{0-sh}(s)$  much larger than the original  $\tilde{R}(s)$ , then  $\Delta I_{0-sh}$  may become strongly downward biased.

It is thus important to assess numerically the dependence of the sampling properties of  $\Delta I_{0-sh}$  on the correlation strength. To address this issue, we next simulated spikes trains as a first order ( $q=1$ ) Markov process with a firing rate  $\rho$  to each stimulus and a given coefficient  $c$  quantifying Pearson correlation of spikes across adjacent time bins (both parameters were taken as constant over time). Reference values for this parameters  $\rho_0$  and  $c_0$  were, similarly to previous simulations, measured (as an average over a post-stimulus time window of 20-40 ms post stimulus onset) from the experimental cortical responses recorded in response to 49 different whisker vibration stimuli by Arabzadeh et al. (2004). In the simulations, we then took the same rate  $\rho_0$  as the real data and we produced a simulated correlation strength  $c$  that was  $f$  times stronger than the real one, i.e.  $c = f c_0$ . Fig. 4 reports how the bias of  $H(\mathcal{R}|\mathcal{S})$  and of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  depends the correlation strength (no bias correction procedure was applied). Fig. 4A shows the results obtained for  $f = 0$ , i.e. absence of correlations. In this case, the number of possible responses is, on average, unchanged by shuffling. Thus  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  have exactly the same bias, and  $I_{sh}$  is unbiased. For  $f = 1$  (i.e. real neuronal correlation value, shown in Fig. 4B), the bias of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  is still very close to that of  $H(\mathcal{R}|\mathcal{S})$ . For correlation four times stronger than that of the real neuron (Fig. 4C) there is some difference between the biases of  $H(\mathcal{R}|\mathcal{S})$  and of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$ . However, with 64 trials per stimulus or more, the bias of the two entropies cancels out and  $I_{sh}$  is unbiased.  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  becomes considerably more biased downward than  $H(\mathcal{R}|\mathcal{S})$  only for correlation eight times stronger than that of the real neuron (Fig. 4D). Thus, we conclude that, although their sampling properties will get worse as the amount of correlation in the data grows, the shuffled estimators  $\Delta I_{0-sh}$  and  $I_{sh}$  can be very useful for the analysis of a wide range of neurophysiological experiments.

## 4.2 The Markov-model decoder

The above section dealt with the use of the simplest decoding model,  $P_0(\mathbf{r}|s)$ . How do the sampling properties of  $I_{LB-q}$  and  $\Delta I_q$  change when using more detailed decoding models such as  $P_q(\mathbf{r}|s)$  with  $q > 0$ ?

To address this issue, we now consider the properties of  $\chi_q(\mathcal{R})$  and  $H_q(\mathcal{R}|\mathcal{S})$  with  $q > 0$ . As  $q$  increases,  $\chi_q(\mathcal{R})$  is expected to become more biased because it depends on  $P_q(\mathbf{r}|s)$  logarithmically. The bias of  $H_q(\mathcal{R}|\mathcal{S})$  is expected to increase even more significantly with  $q$ . In fact it can be decomposed into a sum of stimulus-conditional entropies of the marginal probability distributions of up to  $q + 1$  consecutive time bins (Pola et al., 2003). For this reason, although the bias of  $H_q(\mathcal{R}|\mathcal{S})$  is smaller than that of  $H(\mathcal{R}|\mathcal{S})$ , it grows for larger  $q$  values. The bias of  $I_{LB-q}$  is given by the difference between the biases of  $\chi_q(\mathcal{R})$  and  $H_q(\mathcal{R}|\mathcal{S})$ . Thus, as  $q$  increases,  $I_{LB-q}$  will be more biased.

These expectations are confirmed by the simulations shown in the Fig. 5. These simulations were performed as in Fig. 1, by creating a somatosensory-like, simulated spike train over  $L=10$  time bins responding to 49 different stimuli with realistic firing rates and correlations described by a  $q = 3$  Markov model. A quadratic extrapolation procedure was used to correct for the bias.  $I_{LB-q}$ ,  $\Delta I_q$  and  $\Delta I_{q-sh}$  were computed for  $q = 1, 3$  and  $6$  and plotted in panels A, C, and E respectively. Fig. 5 confirms the intuition that, as the value of  $q$  increases,  $I_{LB-q}$  becomes more biased. The asymptotic value of  $I_{LB-1}$  can be computed well even with as small as 32 trials (Fig. 5A). However, since the actual length of the Markov model used to generate the data was equal to  $q = 3$ , the asymptotic value of  $I_{LB-1}$  is only 0.69 bits, which is less than the true value of the full mutual information carried by the spike train (0.82 bits). This correction reflects the fact that some information would be lost by a decoder neglecting correlations of range longer than one bin. On the contrary,  $I_{LB-3}$  and  $I_{LB-6}$  can reach the correct asymptotic value of information (see Fig. 5C and E). However, they are considerably more biased than  $I_{LB-1}$ .

Figs. 5A, C and E show that the bias of  $\Delta I_q$  has a completely different dependence on  $q$ : the higher  $q$ , the smaller the bias of  $\Delta I_q$ . This is because, as  $q$  increases, the difference  $H(\mathcal{R}|\mathcal{S}) - H_q(\mathcal{R}|\mathcal{S})$  becomes smaller, and thus the overall bias of  $\Delta I_q$  decreases. Similarly to the  $q = 0$  case considered in the previous section, the bias in  $\Delta I_q$  can be further reduced by using a shuffling procedure. In fact it is easy to construct a “ $q$ -shuffling” procedure that generate simulated responses  $\mathbf{r}_{q-sh}$  preserving all marginals and correlations up to  $q + 1$  consecutive time bins, but destroying all the higher length correlations<sup>9</sup>. As in Eq. (23), we

---

<sup>9</sup>The “ $q$ -shuffled” responses  $\mathbf{r}_{q-sh}$  can be constructed as follows. The response in the first time bin is chosen as a random permutation of the first time bins across all trials. Then the response in each bin  $r(k)$  for  $k = 2 \dots L$  is taken randomly without replacement from the subset of trials that satisfy  $r'(k - j) = r(k - j)$  for  $j = 1 \dots m$ , where  $m = \min(k - 1, q)$ . In other words,  $q$ -shuffled bins are concatenated with other taken at random without replacement from the subset of trials that have the same state of up to  $q$  past bins. Naturally, as  $q$  approaches the total length  $L$  of the time window, it could happen that, if data are scarce, the  $q$ -shuffling procedure becomes “trivial” by regenerating always exactly the

can construct from the probabilities of  $q$ -shuffled responses  $\mathbf{r}_{q-sh}$  a “shuffled noise entropy of order  $q$ ”  $H_{q-sh}(\mathcal{R}|\mathcal{S})$ , which converges to the same asymptotic value of  $H_q(\mathcal{R}|\mathcal{S})$  for large numbers of trials, but presents a higher bias for small trial numbers. The introduction of  $H_{q-sh}(\mathcal{R}|\mathcal{S})$ , allows us to define, in analogy to Eq. (26), a shuffled estimator  $\Delta I_{q-sh}$ , as follows:

$$\Delta I_{q-sh} = H_{q-sh}(\mathcal{R}|\mathcal{S}) - H(\mathcal{R}|\mathcal{S}) + H(\mathcal{R}) - \chi_q(\mathcal{R}) \quad (26)$$

$\Delta I_{q-sh}$  must be less biased than  $\Delta I_q$ , because there are substantial bias cancellations between  $H_{q-sh}(\mathcal{R}|\mathcal{S})$  and  $H(\mathcal{R}|\mathcal{S})$ . This expectation is investigated numerically in Fig. 5. Figures 5A, C, and E show the sampling behavior of  $\Delta I_{q-sh}$  for  $q = 1, 3$  and  $6$ , respectively. The bias of  $\Delta I_{q-sh}$  has a negative sign, for the same reasons given above for the case of  $\Delta I_{0-sh}$ . However, for any  $q$  value, the asymptotic value of  $\Delta I_{q-sh}$  for large trial numbers is the same as  $\Delta I_q$ , but they differ significantly for low trial numbers.  $\Delta I_{q-sh}$  is biased downward, and is much less biased than  $\Delta I_q$ .

Figures 5B, D and F consider the behavior of the variances of the information-theoretic quantities as the value of  $q$  increases. The most notable finding is that, like their bias, the variances of both  $\Delta I_q$  and  $\Delta I_{q-sh}$  also decrease with  $q$ . Moreover, for the same reasons presented above for  $q = 0$ , the variance of  $\Delta I_{q-sh}$  is never higher than that of  $\Delta I_q$ . Thus, the decrease of bias of  $\Delta I_{q-sh}$  does not come at the expenses of an increase of variance. This stresses the competitiveness of the shuffling procedures for the estimation of  $\Delta I_q$ .

## 5 Tighter upper and lower information bounds using model selection

We introduced above information theoretic quantities  $\Delta I_q$  that quantify, for any value of memory length  $q$ , the information theoretic cost of using a simplified decoding model described by a Markov model of order  $q$  rather than by the full response probability distribution. Intuition suggests that the shorter the memory needed to decode the neuronal spike trains, the less data are needed to compute its information content. However, we have not explored yet what are the specific advantages offered by *knowing* that, to decode all information, we only need to consider response chains extending over a number  $w$  of time bins which is shorter than the  $L$  time bins making up the time window used to analyze spike trains.

Suppose that we know that the shorter history depth needed to decode all information is  $w$  time bins. This amounts to require that  $\Delta I_q = 0$  for any  $q \geq w$ , and will happen if  $P(s|\mathbf{r})/P_q(s|\mathbf{r})$  is not stimulus modulated for each response and for each  $q \geq w$ . This condition will be met if the stimulus-conditional probability

---

original dataset for each random shuffling. Thus, care should be taken to verify that this is not happen with the data under analysis

of the neuronal responses  $P(\mathbf{r}|s)$  is generated according to a Markov process with length not higher than  $w$  for every stimulus. What are the further advantages offered to us in computing information quantities when we know the actual value of  $w$ ?

If the shortest Markov order needed to decode all the information in the spike train is equal to  $w$  (i.e.  $\Delta I_q = 0$  for any  $q \geq w$ ), then any  $I_{LB-q}$  (with  $q \geq w$ ) will be equal to the total mutual information  $I$ . Therefore this knowledge offers a huge advantage when computing  $I$ . In fact in this case it will be most convenient to compute the true value of  $I$  by using  $I_{LB-w}$ , which will be equal to  $I$  in the asymptotic sampling regime, but it will be much less biased than  $I$  if  $w$  is much shorter than the window length  $L$ <sup>10</sup>.

The knowledge that the shortest Markov order needed to all the information encoded in the spike train the data is equal to  $w$  offers also an advantage in the estimation of the quantities  $\Delta I_q$ , which may still be larger than zero if  $q < w$ . Using the fact that  $\Delta I_w = 0$ , it is easy to show that all  $\Delta I_q$  for  $q < w$  are in this case equal to a simpler and much less biased quantity, called  $\delta I_q$ , and defined as follows<sup>11</sup>:

$$\delta I_q \equiv \Delta I_q - \Delta I_w = H_q(R|S) - H_w(R|S) + \chi_w(R) - \chi_q(R) \quad \text{if } q < w \quad (27)$$

Since  $\Delta I_w = 0$ , it is clear that  $\delta I_q = \Delta I_q$ . However, when estimating the quantities from finite samples,  $\delta I_q$  is much less biased than  $\Delta I_q$ . If  $w < L - 1$ , then term  $H_q(R|S) - H_w(R|S)$  (the dominant term for the bias in Eq. (27)) will have a resulting upward bias that will be smaller than the bias in the corresponding term in the definition of  $\Delta I_q$ , namely  $H_q(R|S) - H(R|S)$ . The smaller  $w$  with respect to  $L$ , the better the sampling advantage in using  $\delta I_q$ .

Along the same lines explained above, we can use the “q-shuffling” procedure to obtain a very weakly downward-biased estimator of  $\delta I_q$ . By computing the noise-entropy difference  $H_q(R|S) - H_w(R|S)$  by using the “q-shuffled” distribution we can define the following quantity:

$$\delta I_{q-sh} = H_{q-sh}(R|S) - H_{w-sh}(R|S) + \chi_w(R) - \chi_q(R) \quad (28)$$

In this case, the difference of entropies  $H_{q-sh}(R|S) - H_{w-sh}(R|S)$  will have only a very small downward bias.  $H_{q-sh}(R|S) - H_{w-sh}(R|S)$  will be less biased than its counterpart  $H_{q-sh}(R|S) - H(R|S)$  in  $\Delta I_q$ . Thus,  $\delta I_{q-sh}$  will be a tighter downward biased estimator than  $\Delta I_{q-sh}$ .

In summary, knowing the true shortest length of the Markov model  $w$  that can decode all information encoded by the spike train is useful because it leads

---

<sup>10</sup>In this case one could in principle use any other  $I_{LB-q}$  (with  $q > w$ ) to compute information  $I$ . However it is more efficient to use  $I_{LB-w}$  because the bias of  $I_{LB-q}$  grows with  $q$ , see Fig. 5.

<sup>11</sup>The quantity  $\delta I_q$  should have a further index  $w$ , which we omit for compactness of notation.



to computing  $\Delta I_q$  through two tight estimators,  $\delta I_q$  and  $\delta I_{q-sh}$ , which bound precisely the information from above and below respectively, and that are tighter and less biased than the corresponding estimators  $\Delta I_q$  and  $\Delta I_{q-sh}$  obtained in the absence of knowledge about  $w$ .

In the next two subsections we shall illustrate better how the knowledge about  $w$  can reduce the sampling bias in a very dramatic way. We will consider two cases. The first case is when we are given some independent knowledge of the value of  $w$ ; the second is when we do not have this knowledge and we have to guess  $w$  from the data.

## 5.1 Using a prior knowledge of the shortest Markov order needed to decode all information

There may be situations in which there is precise knowledge about the memory length  $w$  needed to decode all the information encoded in the spike train. For example, when analyzing the information properties of some simulated data, we often have theoretical insights into the value of  $w$ . Alternatively, for some well studied neural systems there may be data available indicating that there is no correlation between spike times exceeding a certain time separation of  $w$  time bins. In this case, we are entitled to use straightforwardly the quantities  $\delta I_q$  and  $\delta I_{q-sh}$  as bounds to the true value of  $\Delta I_q$ .

In Fig.6A we analyze again the simulated cortical responses to 49 stimuli (generated as in Fig 1-2 over  $L = 10$  time bins) and show the behavior of various information estimators. (Here, a quadratic extrapolation bias corrections was used). As usual,  $I$  and  $\Delta I_0$  are strongly biased upward and require at least  $2^{10}$  trials per stimulus to give estimates close to the asymptotic values. In contrast,  $\Delta I_{0-sh}$ , which is also computed without any assumption on the order  $w$  of the Markov model underlying the real spike train, gives an accurate estimate at around 256 trials per stimulus. What is the effect of making use of the knowledge that the true order  $w$  of the underlying Markov model is equal to 3? This amounts to using  $\delta I_q$  and  $\delta I_{q-sh}$  in Eqs. (27) and (28) with  $w = 3$ . Fig 6A and B shows that  $\delta I_0$ ,  $\delta I_{0-sh}$  provide much tighter and data robust upper- and lower-bound estimations than  $\Delta I_q$ ,  $\Delta I_{q-sh}$ . As expected by the above considerations,  $\delta I_0$  is a much tighter upward biased estimator than the direct use of  $\Delta I_0$ . This is because the knowledge of  $w = 3$  ensures that we compute the functionals in Eq. (27) on a response space with  $2^4$  different responses, a number much smaller than the  $2^{10}$  possible responses characterizing the probability distribution over 10 bins. The downward biased shuffled estimator  $\delta I_{0-sh}$  is even tighter than  $\Delta I_{0-sh}$ , as provides reliable estimates of information with as little as 64 trials per stimulus.

The standard deviation of the estimators  $\delta I_q$  and  $\delta I_{q-sh}$  is considered in Fig. 6B. It is apparent to see that the reduction in bias of  $\delta I_q$  and  $\delta I_{q-sh}$ , with respect to  $\Delta I_q$  and  $\Delta I_{q-sh}$ , is not obtained at the expenses of an increase in

variance.

In Fig. 6C we show the performance of the same estimators obtained with the NSB method bias correction method rather than with the quadratic extrapolation procedure. The comparison with Fig. 6A shows that the direct estimation of  $I$  and  $\Delta I_0$  are less biased than those obtained with the quadratic extrapolation. However, this increase in performance is greatly enhanced when using the  $\delta I_0$  and  $\delta I_{0-sh}$ . The NSB method is in general effective, and  $\Delta I_{0-sh}$  and  $\delta I_{0-sh}$  gave a good estimate of the asymptotic value of  $\Delta I_0$  even down to 128 trials per stimulus. A potential problem with using the NSB method in this context is that, although it performs well, in general does not preserve the sign of the residual bias. Therefore, for very low trial values it could happen that the downward bias quantities  $\Delta I_{0-sh}$  and  $\delta I_{0-sh}$  become occasionally higher than the asymptotic values of  $\Delta I_0$ .

Fig. 6D compares the standard deviations of the estimates obtained with the NSB method. A comparison of the estimates and of the performance of the quadratic extrapolation suggests that the reduction of bias obtained using the NSB method does come at the expenses of a relatively small but appreciable increase in variance, and that this increase is higher for the “shuffled estimators”.

## 5.2 Selecting the order of the decoding model with a non-parametric test

In general, when analyzing real spike trains recorded during a neurophysiological experiment, we do not have a precise prior knowledge of the temporal extent of correlations between spikes. How can we take advantage of the sampling properties offered by  $\delta I_0$  and  $\delta I_{0-sh}$  to the case in which we do not know the value of  $w$  a priori? In this subsection, we address this problem and we show that non-parametric statistical techniques (Efron and Tibshirani, 1993) can be used to compute empirically an *effective* shortest Markov order  $w$  needed to decode all information encoded by the spike train.

The crucial property in the derivation of the expression for  $\delta I_q$  was that  $\Delta I_q = 0$  for all  $q \geq w$ . Therefore we can suggest a simple way to determine effectively a statistical procedure to estimate the value of the memory length  $w$  to be inserted in the definition of  $\delta I_q$ . This value can be defined as the  $w$  representing the the highest value of Markov order above which all  $\Delta I_q$  are zero for any  $q > w$ . In practice, this is not easy to determine from data as individual values of  $\Delta I_q$  that should be zero in the asymptotic large- $N$  limit may be estimated as greater-than-zero just because of statistical fluctuations or because of a residual bias not entirely corrected for (see e.g. Fig. 5). However, the likelihood of an individual outcome of  $\Delta I_q$  to be significantly greater than zero could be established easily with a statistical test evaluating the null hypothesis that  $\Delta I_q$  is zero. Once a particular statistical test has been chosen, the value of  $w$  can be determined as

follows. We start from the largest possible value of  $q$  for which  $\Delta I_q$  can be  $> 0$  (i.e.  $q = L - 2$ ), and we use the statistical test to determine whether  $\Delta I_q$  is significantly  $> 0$ . If this is the case, we cannot discard any correlation of any length and we set  $w = L - 1$ . If instead the null hypothesis that  $\Delta I_{L-2}$  is zero cannot be discarded, we repeat the process by decreasing  $q$  at each step and test again whether  $\Delta I_q$  is significantly different from zero. The parameter  $q$  is then decreased until we find that  $\Delta I_q$  is significantly  $> 0$ . At this point the procedure is stopped, and we take  $w$  as the last value of  $q$  for which the null hypothesis could not be discarded.

A simple way to test the null hypothesis that  $\Delta I_q = 0$  for some  $q$  is to use the following non-parametric “bootstrap” test (Efron and Tibshirani, 1993). We first generated a set of “bootstrapped” data-sets obtained by the “ $q$ -shuffled” procedure describe above by destroying all correlations longer than  $q - 1$  bins. We then use these  $q$ -shuffled responses to compute  $\Delta I_q$ . For each random realization of these shuffled responses, the corresponding  $\Delta I_q$  must be zero. However, because of statistical fluctuations and of errors in fully eliminating the bias, the distribution of  $\Delta I_q$  will peak above zero. We can use this empirically generated distribution to set boundaries on accepting the null hypothesis that the value of  $\Delta I_q$  computed with the original dataset is zero.

In our simulations, we found that the distribution of values of  $\Delta I_q$  computed on the  $q$ -shuffled distribution were approximated by a Gaussian (data not shown). Thus, our statistical criterion for rejection was simply that the actual value of  $\Delta I_q$  computed with the original dataset was higher than 2 standard deviations above the mean of the bootstrapped distribution (the cut-off at 2 standard deviation was chosen because it gave the best estimates of information quantities when applied to the simulated data shown below).

In Fig. 7A-B we show the resulting histograms after applying the bootstrap test to simulated data. The simulations were again the same as in Figs. 1 and 2, simulating a somatosensory neuron with underlying memory length  $w = 3$ , and a response word of length  $L = 10$ . In Fig. 7A, we plotted the distribution of estimated  $w$  values obtained with the bootstrap test when 128 trials per stimulus were available. This distribution is broad and peaked at the correct order of the simulations. As the number of trials per stimulus increases to 256 (Fig. 7B), the test becomes more effective as more of the estimated values of  $w$  are correct. The trend continues as the number of trials increases and the test almost always reports the correct  $w$  for number of trials larger than 512 (results not shown).

In Fig. 8A-B we plotted  $\delta I_0$ , and  $\delta I_{sh-0}$ , computed from Eqs. (27) and (28) with  $w$  determined for each simulation through the test described above (values were corrected with the quadratic extrapolation method). Fig. 8A shows that, even in the case in which  $w$  is not known a priori but must be determined through a statistical test, the lower bound given by  $\delta I_{0-sh}$  has already reached the asymptotic value of  $\Delta I$  for as small as  $2^7$  trials per stimulus. The upper bound,  $\delta I_0$  converges more slowly towards the asymptote, it represents a significant im-

provement with respect to the direct estimation of  $\Delta I_q$ .  $\delta I_{0-sh}$  and  $\delta I_0$  bound the true value of  $\Delta I$  with an error of less than 10% for 256 trials per stimulus, and  $L = 10$ . Thus,  $\delta I_{0-sh}$  and  $\delta I_0$  behave much better than  $I(R, S)$ ,  $\Delta I$ , and  $\Delta I_{0-sh}$ , quantities that do not depend upon the determination of the order of the process.

In Fig. 8C we report the information values estimated using the NSB method. It is apparent from Fig. 8C that both  $\delta I_{0-sh}$  and  $\delta I_0$  perform much better than for the quadratic extrapolation case shown in Fig. 8A. The estimates bound the exact value of  $\Delta I_0$  less than 5% for as low as 128 trials per stimulus<sup>12</sup>.

The above results were obtained using simulated data generated by a Markov process of order  $w = 3$ . To check that this procedure worked also for data generated by higher order Markov processes, we repeated the analysis on synthetic data simulating a somatosensory neuron with different values of Markov order ranging from  $w = 3$  to  $w = 8$ , and the usual response word of length  $L = 10$ . We found results consistent with that of Fig. 7 and 8 (data non shown). In brief, the bootstrap test continued to give estimates of  $w$  peaked around the correct highest value of  $w$  with  $\Delta I_w > 0$  (which in this simulation was constructed to be equal to the  $w$  value used to generate the data). The distributions of values obtained with the bootstrap test were slightly narrower for simulated processes with higher  $w$  values, probably because  $\Delta I_q$  has a smaller variance for higher  $q$  values (see Fig. 5). As the number of trials increased, the information quantities  $\delta I_0$  and  $\delta I_{0-sh}$  always converged (from above and below respectively) to their correct asymptotic value and were less biased than  $\Delta I_0$  and  $\Delta I_{0-sh}$ . Obviously, the lower the Markov order used to generate the process, the bigger was the bias advantage of  $\delta I_0$  and  $\delta I_{0-sh}$  over  $\Delta I_0$  and  $\Delta I_{0-sh}$ .

The above work shows that using simplified model of the Markov family and selecting the most economic decoding model with a statistical test leads to drastic reductions of the bias of the information quantities. However, it is likely that future extensions of this work to include other families of simplified models may lead to even better performance. For example, it is likely that in the presence of serial correlation of adjacent but long interspike intervals, selecting among a class of suffix tree model would perform much better than selecting only among full Markov models, because the depth of history of a suffix tree could be made to depend on when and if the previous spike occurred (Kennel et al., 2005).

---

<sup>12</sup>We remind that the NSB method performed less well than the quadratic extrapolation only when computing low dimensional entropies such as  $H_0(R|S)$  (see Appendix A). Thus, a useful practical consideration is that, when pairing the bootstrap test with the NSB correction method, it is safer to use a quadratic extrapolation to compute entropies such as  $H_0(R|S)$  that appear in Eqs. (27,28) when the statistical test suggests a value of  $w = 0$  as the most likely.

## 6 Analysis of real data

In this section we illustrate the application of the methods developed above to two different datasets of real neuronal recordings from the whisker representation in somatosensory (“barrel”) cortex of rats anaesthetized with urethane.

The first dataset consisted of 21 neural clusters, each recorded from a different electrode in a silicon array made up of a total of 100 electrodes (see Petersen and Diamond (2000) for further details). Spike times from each electrode were determined by a voltage threshold set to a value 2.53.5 times the root mean square voltage. Since it was not possible to sort well isolated units from each channel, spikes from the same recording channel were all considered together as a single neural cluster. It has been estimated that each cluster captured the spikes of  $\approx 2$ -5 neurons; see (Petersen and Diamond, 2000). Neural activity was recorded in response to individual stimulation of one of 9 different whiskers (whisker D2 and its 8 nearest neighbors); individual whiskers were stimulated near their base by a piezoelectric wafer, controlled by a voltage generator. The stimulus was an up-down step function of 80  $\mu\text{m}$  amplitude and 100 msec duration, delivered once per second. The trials per stimulus available were 200-500 for this first dataset.

The time course of the estimates of the information transmitted about stimulus location by spike times from a neuronal cluster in response to instantaneous whisker deflection is reported in Fig. 9A. For the information analysis, we considered the response  $\mathbf{r}$  to be a spike timing code (binarized with temporal resolution  $\Delta t = 5 \text{ ms}$ ) defined in a post-stimulus time window that began at 5-ms post-stimulus and whose end was gradually increased in 5 ms steps up to 55 ms post-stimulus. The 21 clusters in the dataset were analyzed separately and then averaged. The estimation of all the entropy quantities was done using the NSB method (similar conclusions were reached using the quadratic extrapolation; results not shown). The full spike timing mutual information was computed through the estimator  $I_{sh}$  and was found to increase smoothly along the whole time range analyzed. Consistently with the fact that under these conditions neurons typically stop firing after 40-50 ms,  $I_{sh}$  also saturated after  $\approx 40$ -50 ms. We also computed the temporal evolution of  $\Delta I_0$  and  $\Delta I_{0-sh}$ .  $\Delta I_{0-sh}$  was small over all time ranges and accounted for no more than 8% of the total information  $I_{sh}$ .  $\Delta I_0$  was also small but, unlike  $\Delta I_{0-sh}$ , increased markedly for longer time windows. This indicates that  $\Delta I_0$ , unlike  $\Delta I_{0-sh}$ , suffers for sampling problems when the number of bins  $L$  is more than 8-10. The fact that  $\Delta I_{0-sh}$  is a very good estimator of the asymptotic value of  $\Delta I_0$  is confirmed by the computation of the tight upper and lower bounds  $\delta I$  and  $\delta I_{sh}$ . The latter quantities were computed by selecting the  $w$  value required for their computation with the bootstrap test described above. The mean value of  $w$  across channels was  $3.7 \pm 2.8$  (mean  $\pm$  SEM). In Fig. 9A the gray region represents the area enclosed by the upper and lower bounds  $\delta I_0$  and  $\delta I_{0-sh}$ .  $\Delta I_{0-sh}$  remained within the two bounds along the whole time scale considered. In contrast,  $\Delta I_0$  overshoot the data-robust upper

bound  $\delta I_0$  for longer time windows, indicating that it was suffering from sampling problems in that time range. To summarize, this analysis shows that the use of  $I_{sh}$  and  $\Delta I_{0-sh}$  gives precise estimates of the asymptotic values of  $I$  and  $\Delta I_0$  for up to 10 times bins. This an excellent performance considered that the number of trials per stimulus was 200-500. It is much beyond what could be achieved for example with a standard direct method to estimate  $\Delta I_0$ .

The second dataset consisted of 24 recordings of neural clusters, again from rat barrel cortex, obtained with the same type of electrodes and anaesthesia as above. In this case, the stimulation protocol was different. The set of stimuli consisted of 49 different types of sinusoidal whisker vibrations, each defined by a unique combination of amplitude and frequency of vibration, and delivered for 500 ms (see (Arabzadeh et al., 2004) for full details). The trials per stimulus available were 200. As before, we considered the response  $\mathbf{r}$  to be a spike timing code (resolution  $\Delta t = 5$  ms) defined in a post-stimulus time window starting at at 5-ms post-stimulus and gradually increased in 5 ms steps up to 55 ms post-stimulus.

The time courses of the estimates of the information transmitted by spike times from a single neural cluster (averaged over the set of available clusters) about the parameters of whisker vibration are reported in Fig. 9B. In contrast to the previous case, the spike timing Mutual information (computed as  $I_{sh}$ ) continued to grow over time, consistent with the fact that neurons show stimulus-evoked activity for several hundred ms in response to these sinusoidal vibrations (Arabzadeh et al., 2003). As in the previous example, both  $I_{sh}$  and  $\Delta I_{0-sh}$  increased smoothly over time (with no sudden upward jump reflecting a potential bias problem). The accuracy of  $\Delta I_{0-sh}$  in estimating the asymptotic value of  $\Delta I_0$  was substantiated by the fact that it was tightly bound by  $\delta I_0$  and  $\delta I_{0-sh}$  (grey area) up to 55 time bins. ( $\delta I_0$  and  $\delta I_{0-sh}$  were computed by selecting the  $w$  value with the bootstrap test. The estimated mean value of  $w$  across channels was  $6.2 \pm 1.6$ ). In contrast  $\Delta I_0$  overshoot the grey area for long windows, indicating it suffers from an upward bias problem in this time range. A potentially interesting neurophysiological observation is that, unlike in the case of instantaneous whisker deflection,  $\Delta I_{0-sh}$  is now considerable (29% of the total information  $I_{sh}$  at 50 ms post-stimulus), and it grows supralinearly over time. This suggests that correlation may be useful in decoding complex whisker vibrations from cortical neuronal activity, and is consistent with the approximated analytical prediction of (Panzeri and Schultz, 2001c) which suggests that, when neurons keep firing over a sustained period of time,  $\Delta I_0$  is expected to grow at least quadratically as a function of time.

## 7 Discussion

Using information theory to probe the neural code over fine time resolutions, large populations or large time windows remains technically difficult because the analysis of long sequences of spikes requires the collection of large amounts of data to sample the probability of occurrence of each possible spike sequence. Under this condition, information measures suffer from serious bias problems, which impose severe limitations to the range of time scales and population sizes available for analysis (Panzeri and Treves, 1996; Strong et al., 1998).

In this paper we have developed a new procedure to estimate the information carried by spike trains that drastically alleviates its sampling problems. The starting point of the procedure is the observation that, if we break down the mutual information  $I$  into a Kullback-Leibler divergence  $\Delta I_{simp}$  (which bounds the information lost by decoding when ignoring stimulus dependent correlations) and the rest (called  $I_{LB-simp}$ ), then almost all the bias usually comes from  $\Delta I_{simp}$ . The second key observation is that the bias of  $\Delta I_{simp}$  (and, as a consequence, that of the total mutual information) can be drastically reduced (and be made negative) at no increase of variance by an appropriate shuffling procedure. The third key observation is that the bias of  $\Delta I_{simp}$  (and thus that of the mutual information) can be further reduced by the use of a non-parametric test to find the minimal complexity within a class of models of correlation that still permits to compute  $\Delta I_{simp}$  correctly.

From the practical point of view, the overall reduction of the bias is useful because it extends the domain of applicability of information theory. To appreciate how this domain is extended, we note that our numerical examples, based on a response space made of  $2^{10}$  different responses, showed that our techniques eliminated the bias even with number of samples as low as 32 trials per stimulus, a number of the order of the square root of the number of different responses. This is a significant advance with respect over the previous requirement that the number of trials is comparable to the number of different responses. This effectively allows one to *double* the timing precision, the time window length or the population size that can be analyzed. This latter fact is timely because large-scale recording techniques are rapidly becoming available, and because there are theoretical arguments (reviewed in Averbeck et al. (2006)) that suggest that the impact of correlations to the neural code may be particularly important when considering larger populations. The bias reduction techniques open up the possibility of analyzing populations twice as big than those previously considered in information theoretic studies of neural codes. Although this is a significant step forward, it is important to recognize that, due to the “dimensionality curse”, even this advance will not ultimately be enough to the direct analysis of very large populations or of very long temporal sequences at very fine time scales. Alternative approaches that may work better in these more extreme situations include algorithms not relying on binning (Victor and Purpura, 1997; Victor, 2002).

The reduction of the bias problem may also help extending the information analysis to the domain of graded brain signals, such as fMRI or Local Field Potentials (Logothetis, 2003). These graded signals are more difficult to analyze because, unlike spikes, they cannot simply be converted into a binary word sampled with a certain temporal precision. The techniques presented here may lend themselves to the analysis of LFPs/fMRI, by first discretising the graded signals into a number of different levels; then characterising the correlation structure among them; and finally fitting it to a low dimensional stochastic model.

A useful property of the “shuffled” information estimator presented here is that, besides reducing the overall magnitude of the bias when compared to a direct procedure, it makes the bias negative rather than positive. This downward bias property is useful in practical studies of neural codes because a finding of significant extra information in spike timing obtained with this new method will ensure that this additional spike timing information is genuine and not an artefact due to sampling problems.

As pointed out above, an important technical step in the reduction of the mutual information bias was the selection, within a predefined class, of the minimal complexity of the simplified model of correlation that still captures the correct value of  $\Delta I_{simp}$ . This was achieved in practice by introducing a parametric family of Markov models to approximate the correlation structure of the real data, and by using a non-parametric “bootstrap” statistical test to select the order  $q$  of the Markov model that best described the neural response. While we showed that the simple non-parametric test and the simple class of models described here are very effective at reducing the data constraints in information calculations, it is important to note that neither of these steps must necessarily be performed exactly in this form. Particularly interesting families of maximum-entropy simplified correlation models are those considered by Amari (2001) and Schneidman et al. (2006). The model selection can be also performed in different ways, for example through log-likelihood ratio model selection or other types of inference (see e.g. (Cover and Thomas, 1991; Kennel et al., 2005)). An important step of future research is to understand which particular class of models describes neural data more economically, and which statistical model selection technique is more powerful under different circumstances.

In recent years there has been a debate on how best to measure the role of correlations in neural coding (see (Nirenberg and Latham, 2003; Pola et al., 2003; Schneidman et al., 2003) for different points of view). The measures used in this paper is  $\Delta I$ , which was proposed by Nirenberg and Latham (2003). It has an interpretation as an upper bound to the information lost by a decoder that neglects correlation. In this paper  $\Delta I$  was used to break down the information into mildly biased and strongly biased components, and to obtain a more data robust estimator of the total mutual information through this breakdown. The considerable sampling advantages in the computation of the mutual information obtained in this way would be available also to situations when the computation



of  $\Delta I$  is not of interest, either because the only purpose is to quantify mutual information or because other measures of the importance of correlation are used. In fact, the other measures proposed (Pola et al., 2003; Schneidman et al., 2003) all quantify the importance of correlation as differences of certain mutual information quantities; therefore the bias of each such information quantity could be substantially reduced with the techniques presented here. The bias reduction procedure presented here is thus useful to better compute other measures of the importance of correlations on coding.

In summary, the combination of the simplified models and the shuffling methods allowed us to extend the range of applicability of information theory to neuronal responses. Usually, computing accurately information with a straight application of the best bias corrections methods available requires a sample size comparable to the number of possible different neural responses. By applying the same bias correction techniques to the shuffled estimators after a model-selection procedure, it is now possible to estimate accurately information quantities with amounts of data one order of magnitude smaller.

## Acknowledgements

We are grateful to R. Petersen, M. E. Diamond and E. Arabzadeh for many useful discussions and for kindly making available to us the example data used in Fig. 9. Gianni Pola contributed to the very early stages of this work. We are indebted to the anonymous referees for useful insights, particularly on the relation between our work and the maximum entropy principle. This research was supported by the International Human Frontier Science Program (MAM), by Pfizer Global Development (RS), Wellcome Trust 066372, and the Royal Society.

## Appendix A - the NSB method

In this Appendix we sketch some theoretical considerations and we present some additional numerical simulations on the performance of the NSB bias correction method (Nemenman et al., 2004) when entropies are computed from response spaces of different sizes and reflecting different firing rates. While these considerations follow straight from the Nemenman et al. (2004) work, they are helpful to understand why this method is not suited for correcting  $I_{LB-q}$  at low  $q$  values, especially at low firing rates. For simplicity, we focus on correcting the response entropy  $H(\mathcal{R})$ . However, similar considerations hold for the stimulus-conditional entropies.

The NSB method (Nemenman et al., 2004) is rooted in the Bayesian inference approach to estimate the entropy  $H(\mathcal{R})$  (which is a function of a generally unknown underlying probability  $P(\mathbf{r})$ ). The Bayesian approach assumes that there

exists some *a priori* probability density function  $\mathcal{P}_{pr}(P)$  in the continuum space of all the possible probability distributions on the response space  $\mathcal{R}$  with  $R$  elements. The Bayesian estimation of the entropy after the observation of  $\{n(\mathbf{r})\}$  (the experimental number of times  $\{n(\mathbf{r})\}$  in which each response is observed) can be computed as an average of the corresponding entropy over all the possible hypothetical probability distributions weighted by their conditional probability given the data.

Unless we have some other criteria to select a preferred value in the entropy range  $[0, \log_2 R]$ , we would like to have a flat *a priori* distribution of entropy. However, the choice of the prior has a strong impact on the value of  $H^{Bayes}$  unless very many data are observed. Nemenman et al. (2004) have addressed this problem by using a mixture of Dirichlet priors, the latter being defined in terms of a parameter  $\beta$  ranging between 0 and  $\infty$ . Nemenman et al. (2004) have shown that, after fixing  $\beta$  (i.e. after choosing the prior within the family), the Bayesian estimate of the entropy is sharply defined and monotonically dependent on the parameter  $\beta$  (naturally until the number of samples becomes large, in which case the likelihood dominates the estimate). It can be shown that, at fixed  $\beta$ , the variance of entropy estimation before any observation (i.e. when  $n(\mathbf{r}) = 0$ ) scales as  $1/R$  as  $R$  grows, and it is thus small compared to the range of possible entropy values  $[0, \log_2 R]$ . Therefore, the goal of constructing a prior on the space of probability distributions which generates a nearly uniform distribution of entropies can be approximately achieved by an average over the Bayesian estimation over all the one-parameter Dirichlet family of priors labeled by  $\beta$ .

While this procedure is designed to work well for large  $R$ , it may work less well for small values of  $R$ . In fact, if  $R$  is small, then the variance of *a priori* entropy estimation is not small anymore with respect to the range of possible entropies. Thus the result of integration over  $\beta$  will not be flat, but will typically be smaller near the edges 0 and  $\log_2 R$  than in the central region of possible entropies. Thus, the method is likely to give problems in the estimation of low entropy values for low  $R$  values. In particular the NSB method is likely to give problems when estimating the entropy of processes generated with very low firing rates. In this case, since the probability of observing a spike in each bin is low, the entropy value would be much nearer zero than to  $\log_2(R)$ ; since in this case the NSB prior distribution of entropy values is instead higher in the central part of the interval  $[0, \log_2 R]$ , the resulting NSB estimation of a low-firing rate the entropy for low  $R$  may strongly overestimate the entropy, unless very many experimental observations are available. In the latter case, asymptotic bias corrections procedures might work better anyway.

We tested these considerations by applying the NSB method to a homogenous Poisson process. The spike times generated by the Poisson processes were binned using a bin size of  $\Delta t = 5$  ms. In Fig. 10 we report the performance of the NSB by comparing it to the quadratic extrapolation procedure and to the uncorrected measure of the entropy. We tested two different values of  $R$  and two different

firing rates. Panels A and B compare the estimations of the entropy using  $R = 2$  for a firing rate of 2 Hz in A and 40 Hz in B. In both cases the estimation with the NSB method performs worse than the extrapolation procedure of Strong et al. (1998). Panels C and D show the estimators applied to Poisson processes with the same firing rates as in A and B, respectively, and using  $R = 2^6$ . It is apparent now that the NSB performs substantially better than the quadratic extrapolation procedure. We consistently found that, for higher  $R$  values, the NSB method was always more competitive than the quadratic extrapolation (data not shown, but see Fig. 2 for the  $L = 10$  case).

When the NSB estimator is clearly the best estimator (low  $N$  and large  $R$ , see Fig. 10 C-D), it is dominated by the prior. Thus, a potential concern is that this superb performance might be specific to the Poisson process used in the simulation, perhaps because this process is a good match to a distribution within the family of Dirichlet priors. It is conceivable that, for some class of strongly-correlated response distributions, there may not be a good match in the Dirichlet family and thus the superiority of the NSB method in the large- $R$  regime may suffer. To test for this potential problem, we repeated the analysis in Fig. 10 using the same correlated processes used to generate the data in Fig. 1 and 4 (and described in the main text). We found results very consistent with those plotted in Fig. 10 C-D. Although these results cannot rule out completely the above concern, they suggest that the NSB method will perform well in the large  $R$  regime on a wide range of processes with realistic neuronal statistics.

Since the Markov noise entropies  $H_q(\mathcal{R}|\mathcal{S})$  of order  $q$  can be written as sums of entropies defined over  $q$  adjacent time bins, it follows that the NSB method is not suited for correcting  $I_{LB-q}$  at low  $q$  values, especially at low firing rates. It however appears to be an excellent method for correcting the quantities involving entropies defined over long response times, such as e.g.  $I_{LB-q}$  at high  $q$  values.

## Appendix B - A link between our assumptions and the maximum-entropy principle

A direct contact between assumption 2 and the maximum-entropy principle can be made as follows. Consider a one-dimensional family of simplified model  $P_{simp}(\mathbf{r}|s)$  which is obtained from the true  $P(\mathbf{r}|s)$  as a one dimensional trajectory in probability space parametrized by  $\lambda$ :

$$P_{simp}(\mathbf{r}|s) = P(\mathbf{r}|s) + \lambda u(\mathbf{r}|s) \quad , \quad (29)$$

with the “modulator”  $u$  satisfying the normalization  $\sum_{\mathbf{r}} u(\mathbf{r}|s) = 0$ . By computing the zeroes of the derivative of the entropy of  $P_{simp}(\mathbf{r}|s)$  with respect to  $\lambda$ , it is easy to show that the entropy of such  $P_{simp}(\mathbf{r}|s)$  is maximized (as a function

of  $\lambda$ ) when  $\lambda$  is chosen so that:

$$\sum_{\mathbf{r}} u(\mathbf{r}|s) \log_2 (P(\mathbf{r}|s) + \lambda u(\mathbf{r}|s)) = 0 \quad (30)$$

Using Eq.(29), the maximum entropy condition in Eq. (30) can be rewritten as:

$$\sum_{\mathbf{r}} (P(\mathbf{r}|s) - P_{simp}(\mathbf{r}|s)) \log_2 (P_{simp}(\mathbf{r}|s)) = 0 \quad (31)$$

which is exactly the condition (5) requested by our assumption 2. Thus, for any simplified model belonging to the family defined in Eq.(29), the only one that satisfies our assumption 2 is the model with the highest entropy within the family.

This parametric entropy maximization can be related to the construction of the classes of maximum-entropy models developed in Section 3.3, e.g. by considering the extremization of entropy with respect to a large number of modulator functions. Thus, assumption 2 is related to a maximum entropy principle.

## References

- Abeles, M., Bergman, H., Margalit, E. and Vaadia, E. (1993). Spatio-temporal firing patterns in the frontal cortex of behaving monkeys, *J. Neurophysiol.* **70**: 1629–1638.
- Amari, S. (2001). Information geometry on hierarchy of probability distributions, *IEEE Trans. Inform. Theory* **47**: 1701–1711.
- Arabzadeh, E., Panzeri, S. and Diamond, M. E. (2004). Whisker Vibration Information Carried by Rat Barrel Cortex Neurons, *J. Neurosci.* **24**(26): 6011–6020.
- Arabzadeh, E., Petersen, R. S. and Diamond, M. E. (2003). Encoding of Whisker Vibration by Rat Barrel Cortex Neurons: Implications for Texture Discrimination, *J. Neurosci.* **23**(27): 9146–9154.
- Arabzadeh, E., Zorzin, E. and Diamond, M. E. (2005). Neuronal Encoding of Texture in the Whisker Sensory Pathway, *PLOS Biology* **3**(1): 0155–0165.
- Averbeck, B. B., Latham, P. E. and Pouget, A. (2006). Neural correlations , population coding and computation, *Nature Reviews Neuroscience* **7**: 358–366.
- Borst, A. and Theunissen, F. E. (1999). Information theory and neural coding, *Nature Neuroscience* **2**: 947–957.

- Brosch, M., Bauer, R. and Eckhorn, R. (1997). Stimulus dependent modulations of correlated high frequency oscillations in cat visual cortex, *Cerebral Cortex* **7**: 70–76.
- Buracas, G. T., Zador, A. M., DeWeese, M. R. and Albright, T. D. (1998). Efficient discrimination of temporal patterns by motion-sensitive neurons in primate visual cortex, *Neuron* **20**: 959–969.
- Cover, T. M. and Thomas, J. A. (1991). *Elements of information theory*, John Wiley, New York.
- Dan, Y., Alonso, J.-M., Usrey, W. M. and Reid, R. C. (1998). Coding of visual information by precisely correlated spikes in the lateral geniculate nucleus, *Nature Neuroscience* **1**: 501–507.
- de Ruyter van Steveninck, R., Lewen, G., Strong, S., Koberle, R. and Bialek, W. (1997). Reproducibility and variability in neural spike trains, *Science* **21**: 1805–1808.
- DeWeese, M. R., Wehr, M. and Zador, A. M. (2003). Binary spiking in auditory cortex, *J. Neurosci.* **23**: 7940–7949.
- Dimitrov, A. G. and Miller, J. P. (2001). Neural coding and decoding: communication channels and quantization, *Network: Comput. Neural Syst.* **12**: 441–472.
- Efron, B. and Tibshirani, R. (1993). *An Introduction to the Bootstrap*, Chapman and Hall, New York.
- Gawne, T. J. and Richmond, B. J. (1993). How independent are the messages carried by adjacent inferior temporal cortical neurons?, *J. Neurosci.* **13**: 2758–2771.
- Golledge, H. D. R., Panzeri, S., Zheng, F., Pola, G., Scannell, J. W., Gianikopoulos, D. V., Mason, R. J., Tovee, M. J. and Young, M. P. (2003). Correlations, feature binding and population coding in primary visual cortex, *Neuroreport* **14**.
- Gray, C. M., Konig, P., Engel, A. K. and Singer, W. (1989). Oscillatory responses in cat visual cortex exhibit inter-columnar synchronization which reflects global stimulus properties, *Nature* **338**: 334–337.
- Johnson, D. H., Gruner, C. M., Baggerly, K., and Seshagiri, C. (2001). Information-theoretic analysis of neural coding, *J. Comput. Neurosci.* **10**: 47–69.

- Kennel, M. B., Shlens, J., Abarbanel, H. D. I. and Chichilnisky, E. J. (2005). Estimating Entropy Rates with Bayesian Confidence Intervals, *Neural Comp.* **17**(7): 1531–1476.
- Latham, P. E. and Nirenberg, S. (2005). Synergy, Redundancy, and Independence in Population Codes, Revisited, *J. Neurosci.* **25**(21): 5195–5206.
- Logothetis, N. K. (2003). The underpinnings of the bold functional magnetic resonance imaging signal, *J. Neurosci.* **23**: 3963–3971.
- London, M., Schreiber, A., Hausser, M., Larkum, M. E. and Segev, I. (2002). The information efficacy of a synapse, *Nature Neurosci.* **5**: 332–340.
- Merhav, N., Kaplan, G., Lapidot, A. and Shamai Shitz, S. (1994). On information rates for mismatched decoders, *IEEE Trans. Inform. Theory* **40**: 1953–1967.
- Miller, G. A. (1955). Note on the bias of information estimates, *Information Theory in Psychology; Problems and Methods* pp. 95–100.
- Nemenman, I., Bialek, W. and de Ruyter van Steveninck, R. (2004). Entropy and information in neural spike trains: Progress on the sampling problem, *Physical Review E* **69**(5): 056111.
- Nirenberg, S., Carcieri, S. M., Jacobs, A. and Latham, P. E. (2001). Retinal ganglion cells act largely as independent encoders, *Nature* **411**: 698–701.
- Nirenberg, S. and Latham, P. E. (2003). Decoding neuronal spike trains: how important are correlations, *Proc. Natl. Acad. Sci. USA* **100**: 7348–7353.
- Paninski, L. (2003). Estimation of entropy and mutual information, *Neural Computation* **15**: 1191–1253.
- Panzeri, S., Petersen, R. S., Schultz, S. R., Lebedev, M. and Diamond, M. E. (2001b). The role of spike timing in the coding of stimulus location in rat somatosensory cortex, *Neuron* **29**: 769–777.
- Panzeri, S. and Schultz, S. (2001c). A unified approach to the study of temporal, correlational and rate coding, *Neural Computation* **13**: 1311–1349.
- Panzeri, S. and Treves, A. (1996). Analytical estimates of limited sampling biases in different information measures, *Network* **7**: 87–107.
- Petersen, R. S. and Diamond, M. (2000). Spatio-temporal distribution of whisker-evoked activity in rat somatosensory cortex and the coding of stimulus location, *J. Neurosci.* **20**: 6135–6143.

- Pola, G., Petersen, R., Thiele, A., Young, M. P. and Panzeri, S. (2005). Data-robust tight lower bounds to the information carried by spike times of a neural population, *Neural Comp.* **17**: 1962–2005.
- Pola, G., Thiele, A., Hoffmann, K.-P. and Panzeri, S. (2003). An exact method to quantify the information transmitted by different mechanisms of correlational coding, *Network* **14**: 35–60.
- Reich, D. S., Mechler, F., Purpura, K. P. and Victor, J. D. (2000). Interspike intervals, receptive fields, and information encoding in primary visual cortex, *J. Neurosci.* **20**: 1964–1974.
- Rieke, F., Warland, D., de Ruyter van Steveninck, R. R. and Bialek, W. (1996). *Spikes: exploring the neural code*, MIT Press, Cambridge, MA.
- Schneidman, E., Berry, M. J., Segev, R. and Bialek, W. (2006). Weak pairwise correlations imply strongly correlated network states in a neural population, *Nature* **440**: 1007–1012.
- Schneidman, E., Bialek, W. and Berry, Michael J., I. (2003). Synergy, Redundancy, and Independence in Population Codes, *J. Neurosci.* **23**(37): 11539–11553.
- Schultz, S. and Panzeri, S. (2001). Temporal correlations and neural spike train entropy, *Phys. Rev. Lett.* **86**: 5823–5826.
- Shadlen, M. N. and Movshon, J. A. (1999). Synchrony unbound: a critical evaluation of the temporal binding hypothesis, *Neuron* **24**: 67–77.
- Shannon, C. E. (1948). A mathematical theory of communication, *AT&T Bell Labs. Tech. J.* **27**: 379–423.
- Strong, S., Koberle, R., de Ruyter van Steveninck, R. and Bialek, W. (1998). Entropy and information in neural spike trains, *Physical Review Letters* **80**: 197–200.
- Victor, J. D. (2002). Binless strategies for estimation of information from neuronal data, *Physical Review E* **66**: 51903–51918.
- Victor, J. D. and Purpura, K. P. (1997). Metric-space analysis of spike trains: theory, algorithms, and application, *Network* **8**: 127–164.
- von der Malsburg, C. (1999). The what and why of binding: The modelers perspective, *Neuron* **24**: 95–104.

## Figure Captions

**Fig 1: Comparison of the sampling properties of plug-in estimations of different probability functionals.** The plug-in estimators of the probability functionals are plotted as a function of the number of trials per stimulus available. Results were averaged over a number of repetitions of the simulation (decreasing from 200 to 10 as the number of trials per stimulus available increased). We simulated a neuron responding to 49 different stimuli. We considered a time precision of 5 ms and a post-stimulus time window of 50 ms; thus  $L$  was equal to 10. As detailed in the main text, the spike train was simulated with a Markov process with the same mean firing rates and up-to-3<sup>rd</sup>-order marginal probabilities as a real spike train recorded in Arabzadeh et al. (2003) in response to 49 different sinusoidal whisker vibrations. **A)** Average values of plug-in estimators of  $\chi_0(R)$ ,  $H(R)$ ,  $H_0(R)$ ,  $H(R|S)$  and  $H_{0-sh}(R|S)$  obtained without any bias correction. **B)** Average values of plug-in estimators of  $I(R, S)$ ,  $I_{LB-0}$ ,  $\Delta I_0$ , and  $\Delta I_{0-sh}$ . obtained without any bias correction

**Fig 2: Comparison of the sampling properties of estimators computed with the quadratic extrapolation and the NSB bias correction procedure.** Results are plotted as a function of the number of trials per stimulus and were averaged over a number of repetitions of the simulation. The simulated data were obtained exactly as in Figure 1, again considering a post-stimulus window of 50ms discretised into  $L = 10$  bins of size  $\Delta t = 5\text{ms}$ . **A)** and **B): Quadratic extrapolation.** **A)** Averaged estimated values of  $\chi_0(R)$ ,  $H(R)$ ,  $H_0(R)$ ,  $H(R|S)$  and  $H_{0-sh}(R|S)$ . **B)** Averaged estimated values of  $I(R, S)$ ,  $I_{LB-0}$ ,  $\Delta I_0$ , and  $\Delta I_{0-sh}$ . **C)** and **D): NSB estimation.** **C)** Averaged estimated values of  $\chi(R)$ ,  $H(R)$ ,  $H_0(R)$ ,  $H(R|S)$  and  $H_{0-sh}(R|S)$ . **D)** Averaged estimated values of  $I(R, S)$ ,  $I_{LB-0}$ ,  $\Delta I_0$ , and  $\Delta I_{0-sh}$ .

**Fig 3: Variance of shuffling estimators.** The plug-in estimates of the probability functionals were computed (with any bias corrections) from exactly the same simulated neural data as in Figure 1, again with  $L = 10$ . **A)** Scatter plot of  $H_{0-sh}(\mathcal{R}|S)$  vs.  $H(\mathcal{R}|S)$  computed from 100 random realizations of the simulated spike trains, each realization of the simulation consisting of 128 simulated trials per stimulus condition. Each point in the scatterplot represents a data pair taken from the same realization of the numerical simulation. **B)** The variance (across realizations of the simulation) of the estimators of  $I$  and  $I_{sh}$  plotted as function number of trials per stimulus condition. In this case  $I_{sh}$  was estimated using values of  $H_{0-sh}(\mathcal{R}|S)$  from the same realization of the simulation. **C)** Scatter plot of  $H_{0-sh}(\mathcal{R}|S)$  vs.  $H(\mathcal{R}|S)$  computed from 100 random realizations of the simulated spike trains, each realization being made of 128 trials per stimulus condition. Now each point in the scatterplot represents a data pair taken from different randomly-paired realizations of the numerical simulation. **D).** The variance (across realizations of the simulation) of the estimators of  $I$  and  $I_{sh}$  plotted



as function number of trials per stimulus condition. Now  $I_{sh}$  was estimated using values of  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  from different randomly-paired realizations of the numerical simulation.

**Fig 4: The effect of correlations on the performance of the shuffling estimators.** We simulated spikes trains as a first order ( $q=1$ ) Markov process with constant rate  $\rho$  and a given Pearson correlation coefficient  $c$ . The value of  $\rho$  was taken from experimental data (Arabzadeh et al., 2004) as an average over a post-stimulus time window of 40 ms after stimulus onset. The correlation coefficient  $c$  was then adjusted to produce different correlation strengths commensurate with correlation value found in the real data,  $c_0$ . The relationship between the real and the simulated correlation strength was given by a factor  $f$  as  $c = fc_0$ , where  $f$  is a multiplicative factor. Panels **A-D**) plot the bias of  $H(\mathcal{R}|\mathcal{S})$  and  $H_{0-sh}(\mathcal{R}|\mathcal{S})$  as a function of the number of trials per stimulus, for different values of the correlation strength  $f$ . No bias correction was used to compute these entropies.

**Fig 5: Bias and Standard Deviation of estimators making use of simplified Markov models of different orders  $q$ .** We computed  $I_{LB-q}$ ,  $\Delta I_q$  and  $\Delta I_{q-sh}$  for  $q$  1, 3 and 6 as function of the number of trials per stimulus. Results are plotted as a function of the number of trials per stimulus. Panels **A**), **C**) and **E**) report the average of these quantities over a number of repetitions of the simulation, whereas Panels **B**), **D**), **F**) report the standard deviation across simulations. The simulated data were obtained exactly as in Figure 1, again considering a post-stimulus window of 50ms discretised into  $L = 10$  bins of size  $\Delta t = 5$ ms. It is worth noticing that the underlying Markov process to generate the simulated data was of order 3.

**Fig 6: Performance of estimators  $\delta I_q$  and  $\delta I_{q-sh}$  when the Markov order of the underlying model is known in advance.** In addition to  $\delta I_q$ , and  $\delta I_{q-sh}$ , for comparison we also show the average values of  $I(R, S)$ ,  $\Delta I_0$ ,  $\Delta I_{sh}$  as function of the number of trials per stimulus. Panels **A** and **C** report the average of these quantities over a number of repetitions of the simulation, whereas Panels **B** and **D** report the standard deviation across simulations. The simulated data were obtained exactly as in Figure 1, again considering a post-stimulus window of 50ms discretized into  $L = 10$  bins of size  $\Delta t = 5$ ms. The underlying Markov process used to generate the simulated data was of order 3. **A**), Data corrected with the quadratic extrapolation method. **B**), Standard deviations of the data shown in panel **B**. **C**) Estimation based on the method of Nemenman et al. (2004). **D**) Standard deviations of the data shown in panel **C**.

**Fig 7: The bootstrap test to estimate model complexity.** We simulated spike trains as in Fig. 1 using a Markov process of order 3. The post-stimulus time window used in the analysis was of 50ms and was discretized into  $L = 10$

bins with resolution  $\Delta t = 5\text{ms}$ . For each realization (out of a total of 100) of the simulated spike trains we estimated  $w$  using the bootstrap test described in the text. Panels **A)** and **B)** show histograms of the distribution of  $w$  values obtained in this way for 128 and 256 trials per stimulus respectively.

**Fig 8: Performance of estimators  $\delta I_0$  and  $\delta I_{0-sh}$  when the Markov order of the underlying model is selected with the bootstrap test.** In addition to  $\delta I_0$ , and  $\delta I_{0-sh}$ , for comparison we also show the values of  $I(R, S)$ ,  $\Delta I_0$ ,  $\Delta I_{0-sh}$  as function of the number of trials per stimulus. The data correspond to simulated spike trains as in Fig. 1 using a Markov process of order  $w = 3$ . However, the knowledge of the true Markov order  $w$  used to generate the data was not used;  $w$  was instead estimated on each simulation using the bootstrap test described in the text. Panels A and C report the average of these quantities over a number of repetitions of the simulation, whereas Panels B and D report the standard deviation across simulations. **A)** Data corrected with the quadratic extrapolation method. **B)** Standard deviations of the data shown in panel A. **C)** Estimation based on the method of Nemenman et al. (2004) **D)** Standard deviations of the data shown in panel C.

**Fig 9: Analysis of rat somatosensory cortex data.** We show the total mutual information as  $I_{sh}$ ,  $\Delta I_0$  and  $\Delta I_{0-sh}$ . The grey area shows the region delimited by the upper and lower bounds  $\delta I_0$  and  $\delta I_{0-sh}$ . The NSB method was used to correct for finite sampling. **A)** Information about stimulus location conveyed by individual neural clusters recorded in rat somatosensory cortex in response to instantaneous whisker deflections (Petersen and Diamond, 2000). Results reported as average ( $\pm$  SEM) over 21 different clusters. **B)** Information about amplitude and frequency of sinusoidal whisker vibration conveyed by individual neural clusters recorded in rat Somatosensory cortex (Arabzadeh et al., 2004). Results reported as average ( $\pm$  SEM) over 24 different clusters.

**Fig 10: Bias in the NSB method.** Average and std. error of the bias ratio (defined as the bias divided by the true asymptotic value of the entropy) for the NSB, quadratic extrapolation and uncorrected estimators of the entropy. The estimators were applied to homogeneous Poisson process. The generated spike times were binned using a bin size of  $\Delta t = 5\text{ ms}$ . The average and std. error were computed over 1000 realizations of the simulations. Panels **A)** and **B)** correspond to  $R = 2$ , and the rates of the homogeneous Poisson process were 2 Hz in A and 40 Hz in B. Both in A and B the quadratic extrapolation outperforms the NSB method. Panels **C)** and **D)** correspond to  $R = 2^6$ , and the rates of the Poisson processes were respectively the same as in A and B. For larger  $R$  the NSB method performs better than the quadratic extrapolation.

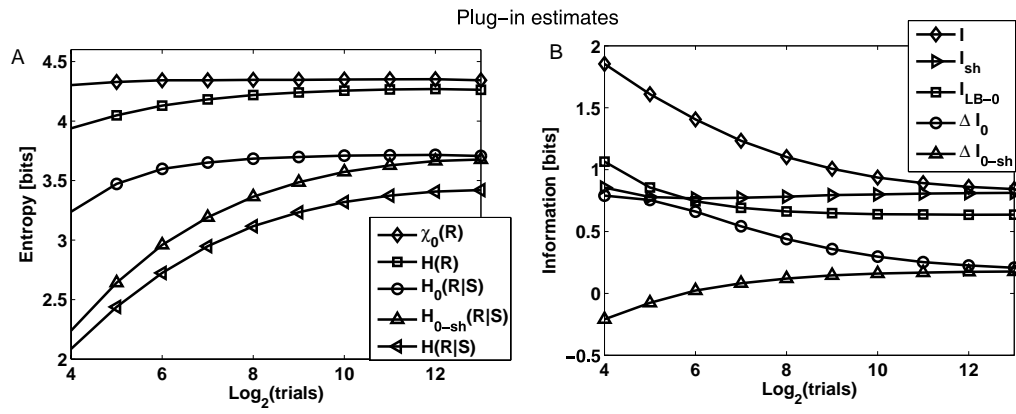


Figure 1: Montemurro *et al*

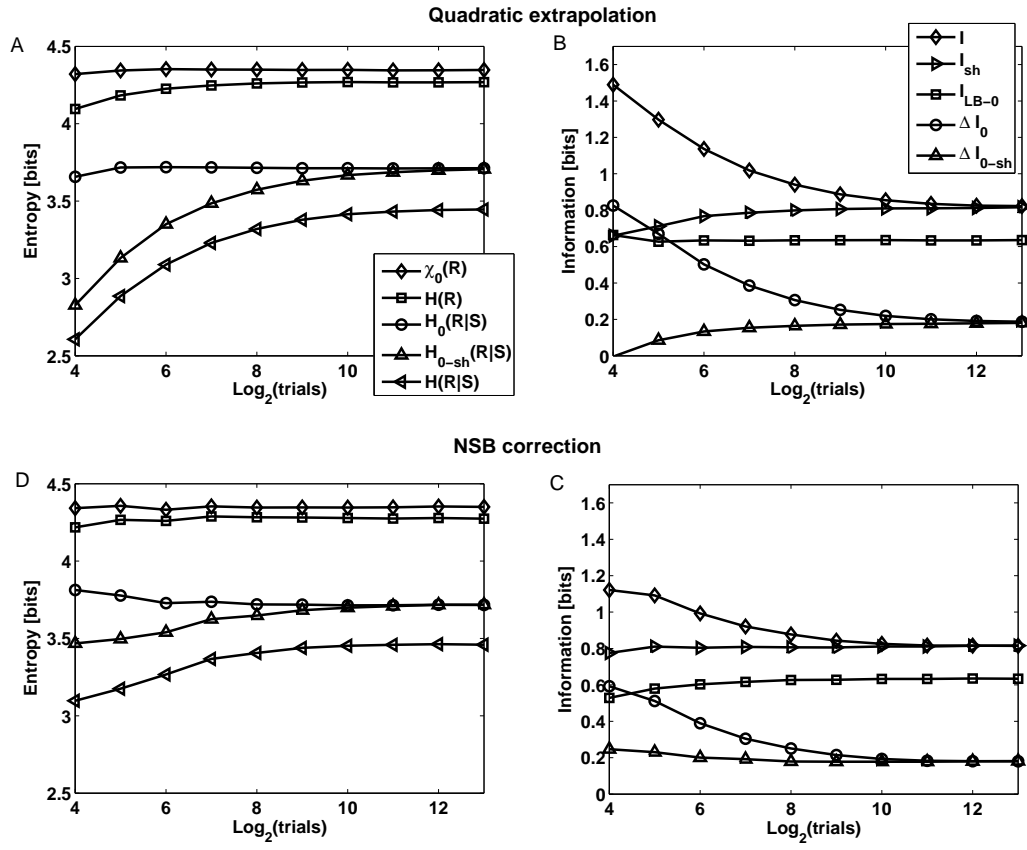


Figure 2: Montemurro *et al*

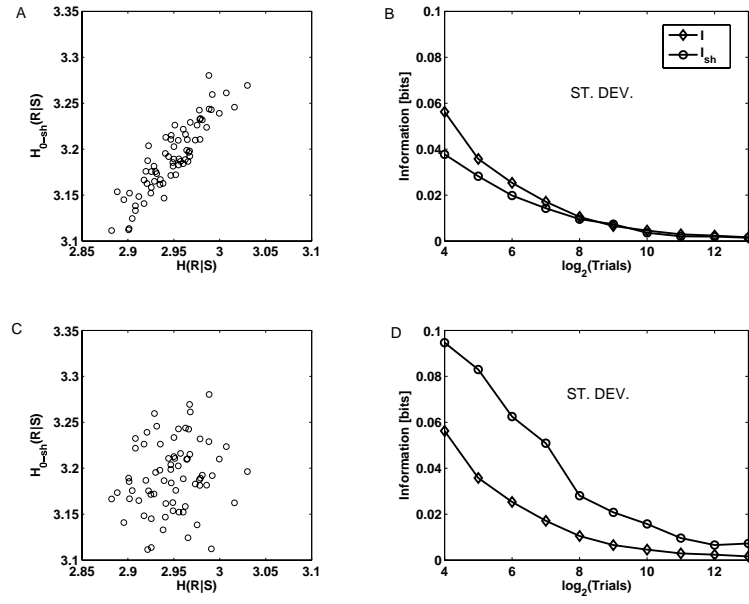


Figure 3: Montemurro *et al*

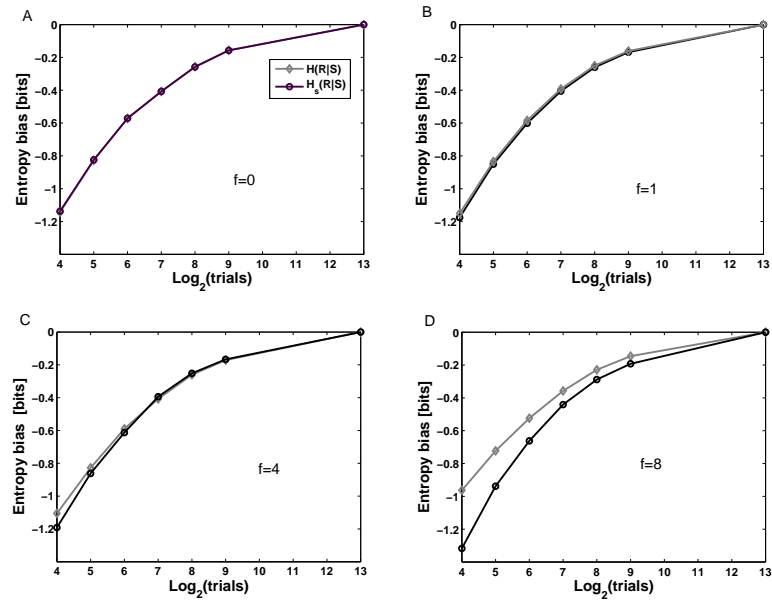


Figure 4: Montemurro *et al*

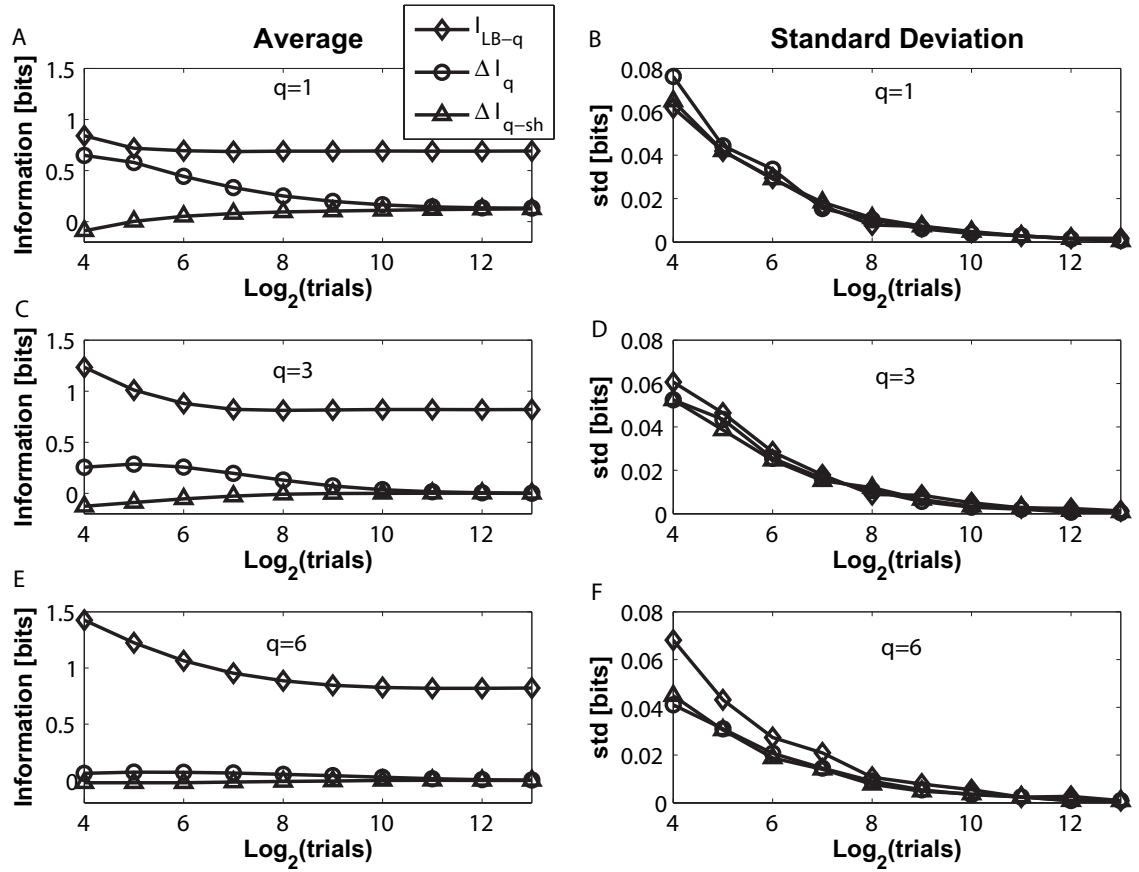


Figure 5: Montemurro *et al*

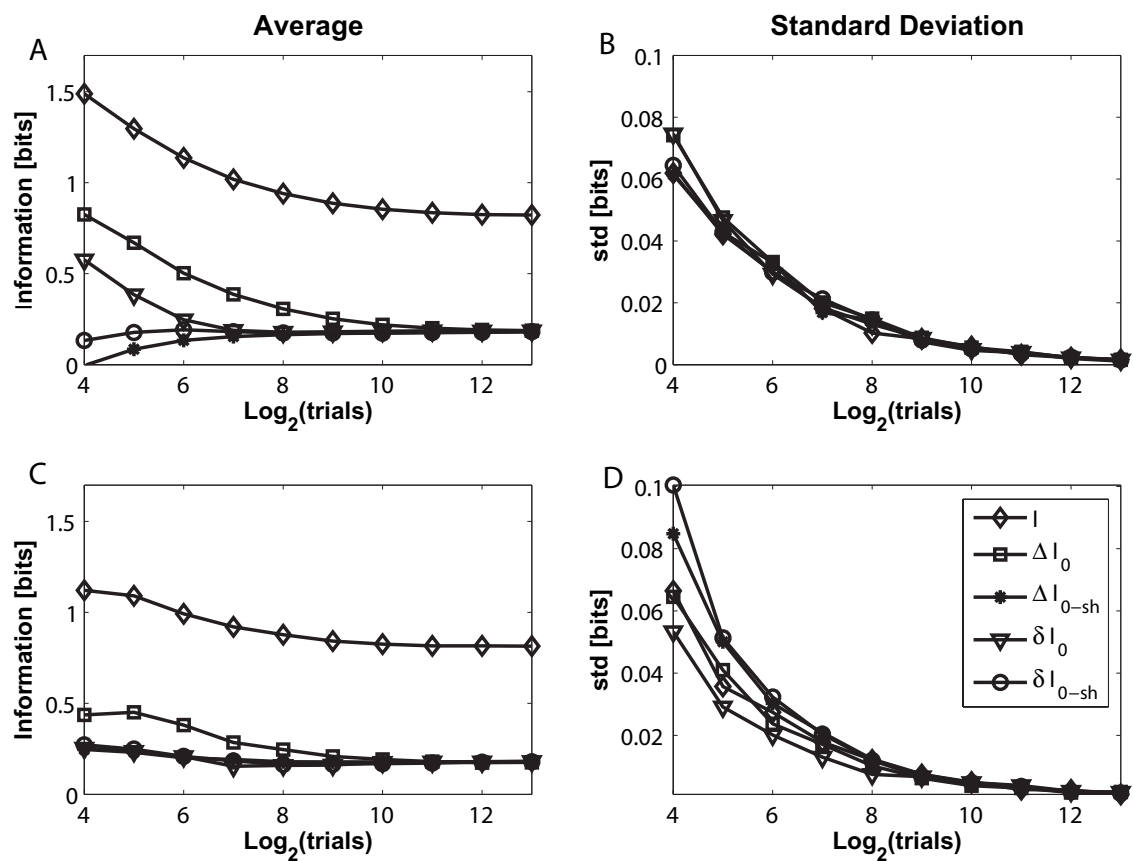


Figure 6: Montemurro *et al*



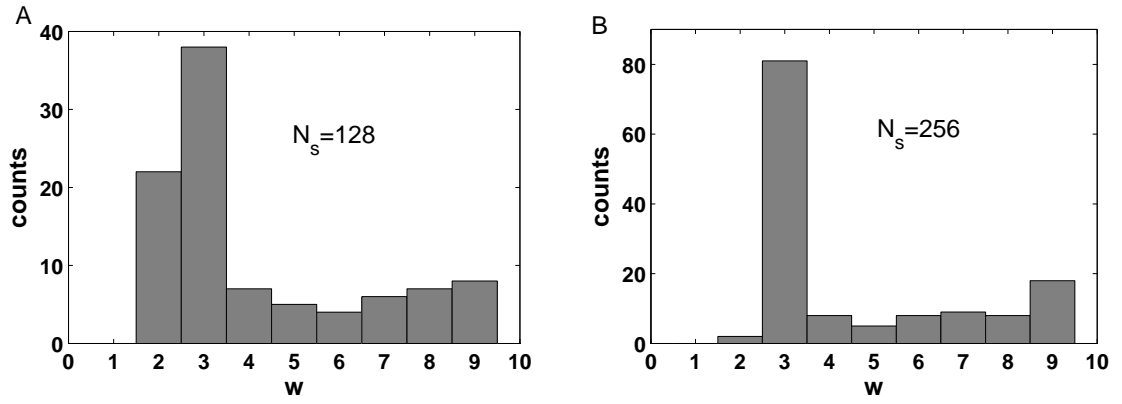


Figure 7: Montemurro *et al*

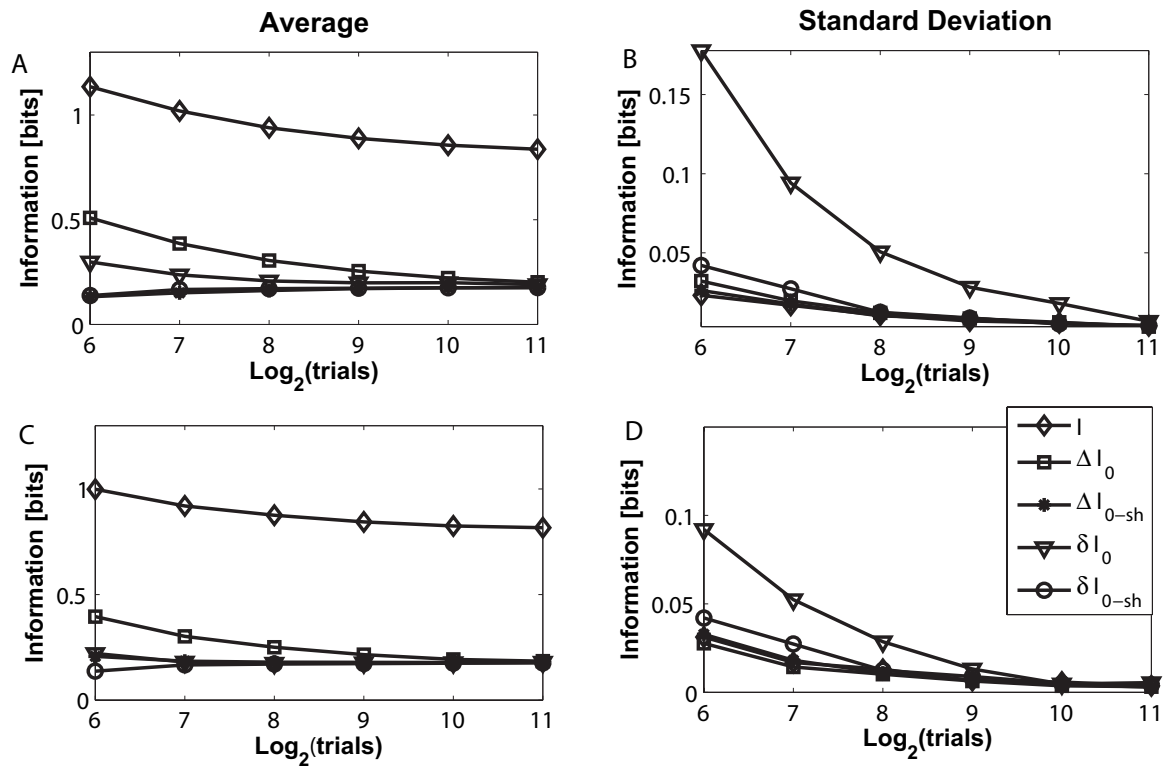


Figure 8: Montemurro *et al*

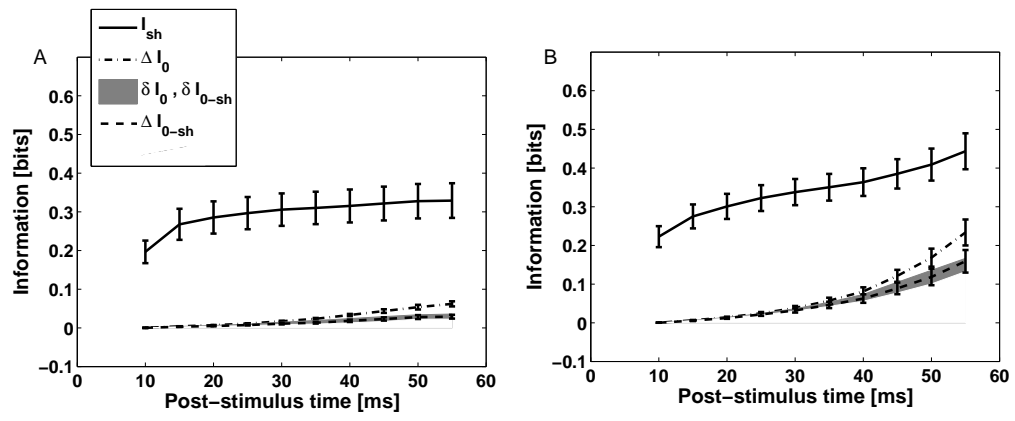


Figure 9:

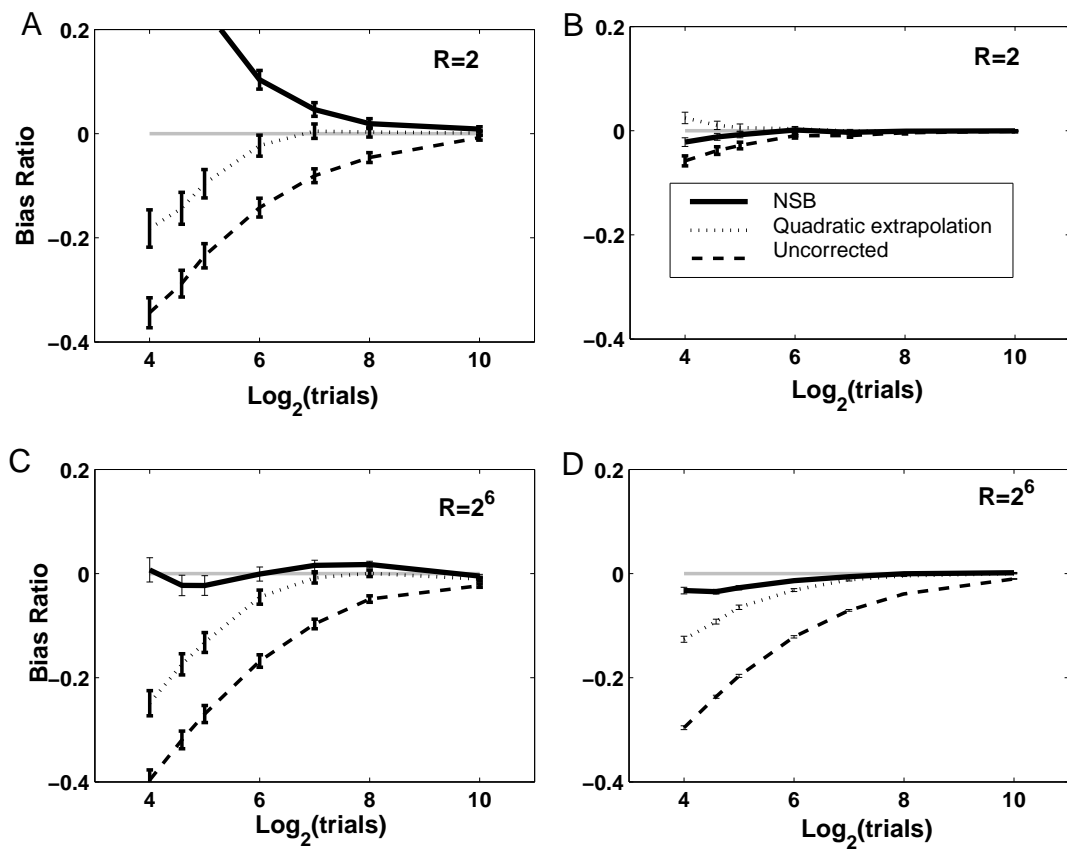


Figure 10: Montemurro *et al*