

1 Review:

2 Characterizing the influence of ‘internal states’ on sensory activity

3 Richard D. Lange* & Ralf M. Haefner*

Brain & Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA

4 March 6, 2017

5 **Abstract**

6 The concept of a tuning curve has been central for our understanding of how the responses of cortical
7 neurons depend on external stimuli. Here, we describe how the influence of unobserved internal
8 variables on sensory responses, in particular their response correlations, can be understood in a
9 similar framework. We review related recent work and discuss its implication for the information
10 contained in sensory responses.

11 **Introduction**

12 A large part of cortical function can be characterized as transforming sensory inputs to behavioral
13 outputs. While this transformation is conceptually unidirectional, the anatomical structures imple-
14 menting it are largely bidirectional [1]. In this review, we will concentrate on discussing progress
15 in understanding feedback signals (FB) that influence neural responses in early visual areas, and
16 their relation to classically feedforward (FF) models of sensory processing. We expect the discussed
17 techniques and computational frameworks to also be useful in understanding neural responses in
18 other sensory or motor areas.

19 The transformation of sensory inputs into behavioral outputs can be understood on at least
20 two levels. On the first one, we would like to understand how a stimulus, S , influences neural
21 responses, \mathbf{r} , and how these responses influence behavior, B . On a more abstract level, one can try
22 to understand how a stimulus affects internal states, I , and how these internal states are linked to
23 behavior (Figure 1). While those two levels are clearly linked, the representation of abstract internal
24 states is generally unclear, and it is in principle only possible to directly observe the corresponding
25 neural responses. Internal states are only accessible by their covariability with one or more of
26 the observable quantities, S , B , and \mathbf{r} . In general, all of these quantities vary over time [2], and
27 observations of their joint or conditional probability distributions can tell us about different aspects
28 of brain function (Figure 1).¹

29 $p(S, \mathbf{r})$: Characterizes the stimulus-dependence of neural responses, giving rise to tuning func-
30 tions and receptive fields.

*{rlange,rhaefne2}@ur.rochester.edu

¹We use joint distributions here because they allow us to characterize relationships even when one of the variables is not directly observed.

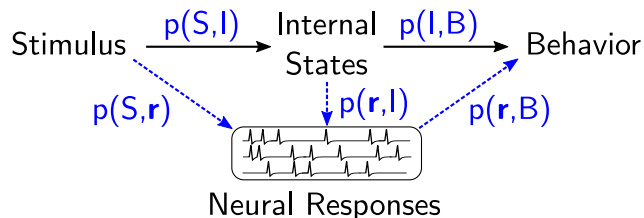


Figure 1: **a:** A computational-level description of the brain may invoke abstract *Internal States* (I) that govern behavior (B) and are influenced by stimuli (S). On a neurophysiological level, we also seek a description on the level of *responses* (**r**) of populations of neurons. It is often useful to mix levels of abstraction, for example modeling the effect of attention (an abstract state) on neural responses.

31 $p(\mathbf{r}, B)$: Usually called ‘decision-related’ or ‘choice-related’ activity in the context of decision-
 32 making tasks [3]. When **r** is measured in motor cortex, it gives rise to tuning curves
 33 with respect to motor outputs [4, 5].

34 $p(\mathbf{r}, I)$: Depending on the nature of *I*, this can alternatively be interpreted as the influence of
 35 cortical states on neural responses (e.g. anesthesia, attentional state, motor activity)
 36 [6, 7, 8, 9^{••}, 10, 11^{••}, 12, 13], or their neural representation (e.g. for beliefs about the
 37 outside world) [14, 15, 16[•], 17[•]].

38 $p(S, I)$: Not observable directly, but may be induced by experimental design, for instance by
 39 training a subject to allocate attention following an external cue, or when measuring
 40 neural correlates of reward and value.

41 $p(I, B)$: Captures differences in behavior depending on internal state, e.g. performance differ-
 42 ences between attentional states, or different choices depending on percept. Again,
 43 not directly observable, but usually implicitly assumed to be under the control of the
 44 experimenter.

45 We will first describe the implications of extending the concept of a tuning function to charac-
 46 terizing the influence of internal states on neural responses. We will then describe recent efforts to
 47 distinguish between feedforward and feedback influences on neural responses, and the implications
 48 for the information represented by sensory neurons.

49 1 Characterization of the influence of internal states

Box 1: Tuning and Covariability

Just as the dependence of a neuron’s firing rate on an *external* parameter yields a tuning curve, a similar dependence on some *internal* state may be thought of as “tuning” to that internal state, though the value of that state may not be precisely known [2]. Taken together, we can write a population’s response as a vector-valued ‘tuning’ function of the external parameters and internal states on which it depends[†]:

$$\mathbf{r} = \mathbf{f}(\underbrace{a, b, \dots}_{\text{external}}, \underbrace{x, y, \dots}_{\text{internal}}) + \text{noise} \quad (1)$$

Taking the linear approximation of \mathbf{f} (Figure 2), the covariance between two neurons’ responses takes the form of a sum of outer products [18[•], 19^{••}, 16[•]]:

$$\mathbf{C} \approx \mathbf{C}^0 + \underbrace{\sum_{v \in \text{vars}} \frac{d\mathbf{f}}{dv} \frac{d\mathbf{f}^\top}{dv} \text{var}(v)}_{\text{independent terms}} + \underbrace{\sum_{u, v \in \text{vars}} \left(\frac{d\mathbf{f}}{du} \frac{d\mathbf{f}^\top}{dv} + \frac{d\mathbf{f}}{dv} \frac{d\mathbf{f}^\top}{du} \right) \text{cov}(u, v)}_{\text{interaction terms}} \quad (2)$$

where \mathbf{C}^0 is ‘intrinsic’ covariance (i.e. the covariance of the noise) and $u, v \in \text{vars}$ represents any of the internal or external variables that f depends on.

Equation (2) makes explicit how shared sources of variability introduce structured covariability into a population. Importantly, both variation in *internal* states and variation in *external* stimulus parameters affect a population’s covariance analogously. By analyzing a population’s covariance structure, then, we can learn about both stimuli and internal states modulating neural responses.

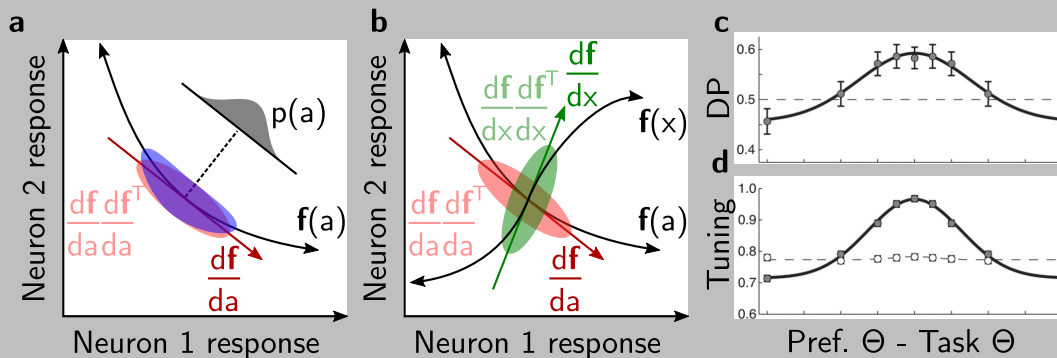


Figure 2: **a**: The population tuning of two neurons to some parameter, $\mathbf{f}(a)$, entails response co-variability along \mathbf{f} (blue) [18[•]], often called ‘signal correlations’ when a reflect the stimulus, and ‘noise correlations’ when it reflects uncontrolled internal variability [20, 21]. **b**: Variability in multiple parameters adds (eq. 2). **c-d**: Population responses of [22] replotted to illustrate conceptual similarity between tuning with respect to outside stimulus (d), and tuning with respect to behavior (c). Here, monkeys performed a motion-based detection task while the authors recorded from macaques’ MT neurons. **c**: Dependence of a neuron’s ‘detect probability’ (DP) on its preferred motion direction relative to the task direction measured. DP is closely related to the difference in the mean responses associated with the two choices, $\Delta\mathbf{f}/\Delta\text{choice}$ [23]. Hence, if DPs reflect FB influences from a decision-making area[24^{••}, 25], then this curve is interpretable as *tuning to the decision state*. **d**: Normalized difference in response associated with the two task-relevant stimuli, $\Delta\mathbf{f}/\Delta\text{stim}$, showing the relationship between tuning to internal and to external parameters.

[†]Correlations between I and \mathbf{r} considered here may reflect causation from I to \mathbf{r} , or common input to both. Equation (2) holds as long as I is not purely a function of \mathbf{r} , in which case variability in I will not affect variability in \mathbf{r} .

50 **Stimulus-dependence of responses $p(S, \mathbf{r})$:** Observing the full joint distribution $p(S, \mathbf{r})$ is
51 rarely possible even for a single neuron due to the high dimensionality of S , at least in the visual
52 domain. As a result, experimentalists have often concentrated on probing this relation along a
53 single dimension while keeping everything else about the stimulus constant. For example, varying
54 the orientation of a grating pattern yields the orientation tuning curve of each recorded neuron.
55 In general, a neuron may be tuned to multiple independent parameters. For example, V1 neurons
56 are often characterized jointly by their tuning to the orientation and to the spatial-frequency of a
57 grating stimulus.

58 **Behavior-related responses $p(\mathbf{r}, B)$:** Taking the same approach as above, we can use concurrent
59 observations of neural responses and behavioral outputs to infer their relationship. Traditionally,
60 this has been done by computing ‘choice probability’ and ‘detect probability’ [26, 25, 27] – measures
61 for how neural responses differ depending on behavioral response. A closely related but simpler
62 measure is the dependence of the *mean* neural response on the behavioral response. For continuous
63 behavioral outputs such as eye or limb movements [28, 29], such a measure is equivalent to the
64 classic concept of a tuning curve. For binary responses as common in perceptual decision-making
65 tasks (reviewed in [25]), it is closely related to choice and detect probabilities under reasonable
66 assumptions [30].

67 In a purely feedforward, unidirectional model, the dependence of neural responses on behavior
68 would reflect the decoding strategy of the brain, i.e. how the neural responses inform behavior.
69 However, in a bidirectional model including internal variables that influence both responses and
70 outputs, this dependence is best thought of as a covariance between responses and outputs, and
71 does not imply only feedforward causation (Box 1).

72 **Internal state-related responses $p(\mathbf{r}, I)$:** Since a subject’s internal state is not observable
73 directly, the traditional approach to measuring its impact on neural responses has been to exploit
74 its relationship to observed quantities like stimulus and behavior. For example, training an animal
75 in an attention task induces a relationship between an aspect of the internal state (‘attentional
76 state’) and an external sensory signal (‘attention cue’) [2], where improved behavioral performance
77 is taken as an indicator for successful manipulation of attentional state [9••]

78 There are two limitations to this approach. (1) there usually remains significant uncertainty
79 about the relationship between observables and internal states. For instance, even though an
80 attention cue may only have two states, the internal variables determining the animal’s attentional
81 state is likely non-binary and more variable. And (2), the internal state is much higher-dimensional
82 than the observations. For instance, the attentional state is likely high-dimensional determining
83 location(s), feature(s), etc. that the animal is attending to [9••, 31, 19••] (Box 2). Equally, the
84 internal percept is surely richer than the binary output in a perceptual decision-making task. Both
85 problems severely limit the information that can be inferred about internal states and their influence
86 on neural responses from their dependence on the preceding stimulus or their dependence on the
87 behavioral report.

88 **What covariability reveals about internal states:** The above approaches can all be consid-
89 ered *conditional*: they rely on measuring each neuron’s response conditioned on other quantities,
90 e.g. conditioned on choice or on the attentional cue. One way to overcome the described limitations
91 is by considering the *full* response distribution, i.e. across all choices or cues. If some component,
92 x , of the subject’s internal state varies uncontrolled from trial to trial, the population response \mathbf{r}
93 moves together according to its ‘tuning’ to x (Box 1). This means that variability in an internal

94 state leaves a signature in the statistical structure of the neural population response. Furthermore,
95 different internal states all leave their own signatures that will superpose (Box 2). By recording
96 from sufficiently many neurons concurrently, and by using sufficiently powerful statistical inference
97 techniques, both of which are under current development [32[•], 33, 34], we can infer *both* the internal
98 states influencing neural responses, and how those responses depend on them [9^{••}, 31, 35^{••}]. This
99 will also allow estimation of the internal state as a function of time [11^{••}] and on individual trials
100 [9^{••}]. Furthermore, with nonlinear dimensionality reduction methods one can infer curvature in \mathbf{f}
101 as well, going beyond the simple but instructive picture of linear tuning presented above [32[•]].

102 **Experimental data:** A number of papers can be interpreted in this framework as having mea-
103 sured consequences of tuning curves with respect to internal variables by correlating neural re-
104 sponses with behavior [22, 36, 37, 38^{••}, 9^{••}, 31, 11^{••}, 35^{••}, 39^{••}], four of which we will highlight
105 here (also see Fig. 2 in Box 1, and Box 2).

106 Having trained a monkey on a cued change-detection task, [9^{••}] computed the line connecting
107 the mean neural response in attention trials with the mean neural response in the attend-away trials.
108 This is equivalent to computing \mathbf{f}' with respect to attention. The authors overcome the challenge
109 that attentional state is not precisely known by binarizing it ('attend left' or 'attend right') and
110 discarding false negative trials, where attention is assumed to not match the cued location. They
111 found that their trial-by-trial estimates of attentional state predicted performance on individual
112 trials, and that it varied substantially across trials even within the same experimental condition.

113 [11^{••}] regressed a scalar latent factor influencing neural activity in primate V1 both during
114 anesthesia and while awake. The shared, time-varying gain term accounted for a large part of the
115 population's variability as well as their correlations during the anesthetized, but not during the
116 awake state. This is an example of inferring an internal state which is clearly important but for
117 which no normative account exists at present.

118 [35^{••}] explicitly modeled the trial-by-trial changes in V4 spiking activity during a change-
119 detection task as a combination of a stimulus-driven component and shared time-varying feedback
120 terms that were either known (the attentional cue) or fit to the data. The inferred values of the
121 fitted internal states were then correlated per-trial with spiking statistics ($p(I, \mathbf{r})$) and the animal's
122 behavior ($p(I, B)$), and were found to explain a large amount of variance in both.

123 At least two groups have measured the effect of changing the task context on noise correlations
124 while keeping the stimulus distribution the same [38^{••}, 39^{••}]. [38^{••}] find, in a change detection
125 task, correlations that are well-described as the effect of only two varying attentional states (Box
126 2). [39^{••}] measured the noise correlation matrix in V1 while the monkey performed a coarse
127 discrimination task. Interestingly, they found that the feedback component of correlations could
128 largely be explained by variability in only a single state. Furthermore, [39^{••}] were able to isolate
129 and quantify the task-dependent component of both choice probabilities and noise correlations,
130 finding a significant contribution of task-related covariability to both.

Box 2: Case Study: [38**]

In Box 1 we argued that response correlations (or more precisely, covariances) may be usefully understood as the sum of ‘intrinsic’ covariability with outer products of the population’s ‘tuning’ to each independently varying parameter. [38**] used essentially this idea, which we replicate and make more explicit here. The authors measured pairwise noise correlations between MT neurons in a motion discrimination task, where the discriminated directions changed from trial to trial. Recording from pairs of neuron at the same time, switching tasks implied that firing by the two neurons either supported the same decision, or opposite decisions. Their results showed noise correlations change depending on whether the neurons support the same choice or different choices, implying an influence of the internal state ‘task context’ on MT responses (Figure 3, bottom left).

Following [38**], we assume two independently varying internal states, x and y : one for alternating ‘attention’ between the discriminated directions, and one for fluctuating ‘attention’ to both of the task-relevant motion directions concurrently [19**].

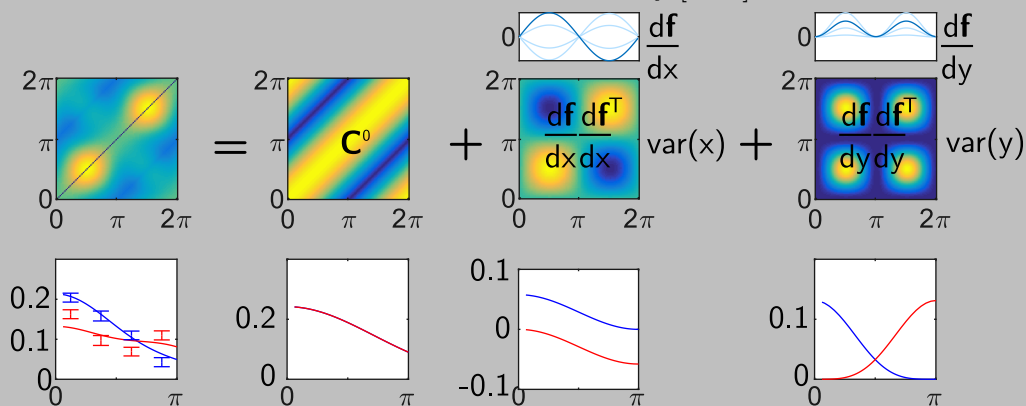


Figure 3: **Middle row:** a full simulated correlation matrix, plotted as a function of each neuron’s preferred motion direction, is derived from the weighted sum of independent variability (not shown), limited-range covariance, and independent outer-product terms (Equation 2). **Bottom row:** [38**] measured pairwise correlations as a function of the difference in the pair’s preferred directions, ranging from 0 to 180°. Data replotted as error bars and compared to model (lines). Pairs where neurons support the same choice are in blue, and pairs supporting different choices are in red. Both, the initial separation between the conditions as well as the crossover at larger differences in preferred direction are captured by the above two components of attention.

131 2 Distinguishing feedforward from feedback influences

132 Because the external stimulus is under the direct control of the experimenter, it is possible to
 133 establish a causal feedforward relationship between external inputs and the responses of sensory
 134 neurons. In order to establish causality in the relationships between neural responses and behavior,
 135 or neural responses and internal states, it is likewise necessary to experimentally manipulate either
 136 of them. Two non-causal approaches to distinguish between feedforward and feedback signals have
 137 been (1) to compare the observations to the predictions from a purely feedforward model and ascribe
 138 any residuals to feedback influences, and (2) to compare neural responses for different internal states
 139 under the *assumption* that those states are represented by neurons downstream of those that are

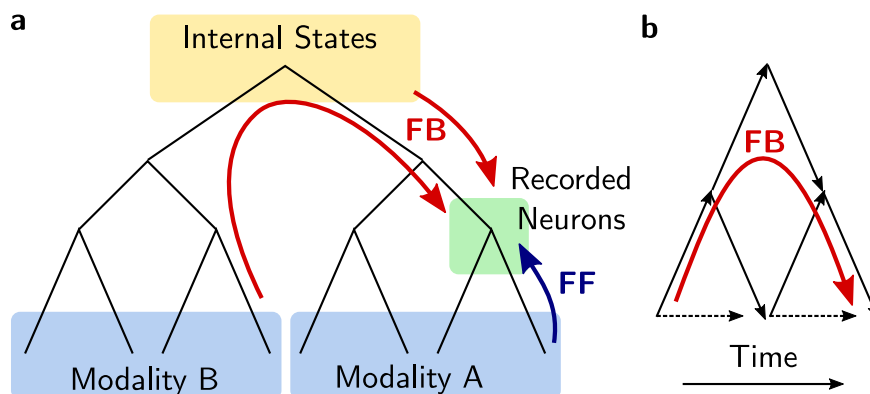


Figure 4: Identification of feedback signals by dependency of neural response on stimuli not within the presumed feedforward pathway. **a:** Dependency of neural responses in other sensory modalities, or outside the receptive field. **b:** Dependency on external inputs at different times, sufficiently long in the past such that ‘memory’ about them within the feedforward pathway can be excluded.

140 being observed (Figure 4). Note that we refer here only to the functional definition of FB, and
141 that FB signals from downstream neurons may themselves not involve just pure cortico-cortico FB
142 pathways but also anatomically FF pathways [40], for instance via a thalamocortical loop [41, 42].

143 [24••] took the first of these approaches. In monkeys performing a disparity discrimination
144 task, they compared the relationships between stimulus and choice to the relationship between
145 sensory area V2 neural responses and choice. They found that the stimulus was most strongly
146 correlated with the choice at the beginning of a trial, while the neural responses were most strongly
147 correlated with the choice towards the end of the trial. Assuming that correlations between the
148 neural responses do not change over the course of a trial one would expect the time courses of
149 stimulus–choice, and response–choice correlations to agree, contrary to what is observed [43]. The
150 authors therefore concluded that the neural responses they observed must be partly due to non-
151 feedforward influences, most notably feedback from downstream decision-related neurons. Since it
152 is conceivable that either feedforward inputs or recurrent connectivity among the sensory neurons
153 could lead to increasing response correlations, answering this question definitively requires direct
154 observations of those correlations.

155 Most studies taking the second approach manipulate the attentional state of the subject, e.g.
156 by presenting a pre-trial cue, and have been extensively reviewed elsewhere [44]. An alternative
157 approach is to analyze the co-variability of the neural responses and determine whether they contain
158 task-dependent components. Under the critical assumption that task-learning, or task-switching,
159 does not alter the feedforward connectivity, or recurrent connectivity, one can then conclude that
160 any such task-dependent component must be due to feedback signals. The underlying assumption
161 is more likely to be true for early sensory areas than higher levels of sensory processing where
162 learning-induced plasticity has observed to be stronger [45]. The assumption is also more likely
163 to be true for task-switches across shorter time scales, e.g. from trial to trial (seconds) rather
164 than days or weeks over which relevant changes of the feedforward pathway become more likely.
165 Both [38••] (see Box 2) and Bondy & Cumming [39••] have taken this approach comparing changes
166 across trials in MT responses, and changes separated by several days in V1 responses, respectively.

167 Even though [35••] did not explicitly vary the task, a similar logic can be applied: Since the
168 main common factors driving trial-to-trial variability impact individual neurons depending on their
169 *task-dependent* tuning properties, they have to be different if the task is changed, and hence they

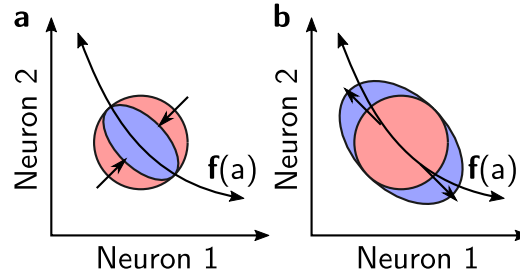


Figure 5: $d'd'$ -correlations ($d' \equiv f'/\sigma$) can arise in different ways, with differing implications for sensory information which is determined by covariance. The blue response distributions have the same correlation but different covariances. **a**: No change in sensory information if $d'd'$ -correlation is induced by suppression of variability orthogonal to \mathbf{f}' . This case implies lower response variability for both neurons. **b**: Decrease in sensory information if $d'd'$ -correlations is induced by increased variability along \mathbf{f}' , e.g. by variability in an internal state whose tuning is aligned with \mathbf{f}' , *assuming that this internal state contains no information about a* .

170 must be due to feedback signals.

171 In this framework, training on a task induces a task-specific association or relationship between
 172 internal states like choice and sensory responses. If those internal states are represented by down-
 173 stream neurons, and if this relationship is purely feedforward in nature, then neural responses will
 174 affect those states, and in turn observable behavioral outputs, but not the other way around. If,
 175 however, those internal states affect sensory responses, then sensory responses will differ depending
 176 on the task.

177 3 Implication for the information in sensory responses

178 Fisher information is often used to quantify the amount of information contained in sensory pop-
 179 ulation responses. It is important to note that this information is computed not with respect to
 180 the entire input signal, S , but with respect to an experimenter-defined aspect, a , of it (e.g. with
 181 respect to orientation instead of the entire retinal image). In general, response variability that is
 182 indistinguishable from the variability along $\frac{df}{da}$ that would be induced by variability in the relevant
 183 aspect of the external signal, limits information with respect to that aspect [20, 18•]. Whether
 184 the response variability induced by variability in internal states increases or decreases Fisher in-
 185 formation depends principally on two questions: first, whether the internal signal itself contains
 186 information about a , and second, whether the way it influences sensory responses is aligned with $\frac{df}{da}$
 187 [19••]. For instance, in the context of a perceptual experiment, an internal variable may represent
 188 the subject's expectations about the stimulus in the upcoming trial. If those expectations bias the
 189 sensory responses in an analogous fashion to actual sensory information, trial-to-trial variability in
 190 those expectations will induce sensory response variability in the $\frac{df}{da}$ direction [16•]. If those expec-
 191 tations are based on superstition (e.g. if trials are randomized), then this internal-state-induced
 192 variability will decrease information. However, if they reflect correctly learned serial dependencies
 193 between trials, then they will increase information.

194 Consider the result that attention reduces covariability in the $\mathbf{f}'\mathbf{f}'^T$ -direction where the deriva-
 195 tive is taken with respect to the task-relevant stimulus [35••]. Since trials were randomized in the
 196 underlying experiment, an internal variable whose value depended on previous stimulus presenta-
 197 tions could not contain any information about the upcoming stimulus. Reducing sensory variability

198 with respect to such a variable would therefore *increase* information in the sensory representation
199 in line with the behavioral improvement experimentally observed [35••].

200 Importantly, information is limited not by correlations themselves but by covariability, the
201 product of correlation and individual variability. Consider Figure 5 which illustrates two different
202 ways in which correlations in the \mathbf{f}' -direction can emerge. Either (panel c), by *increasing* variability
203 that is indistinguishable from stimulus-induced variability and hence reduces information. Or (panel
204 a), by *decreasing* variability in orthogonal directions that leave variability along \mathbf{f}' and, hence,
205 information unchanged [19••].

206 4 Conclusion

207 The influence of internal states on sensory responses can be characterized by tuning functions
208 in analogy to the influence of the external stimulus on neural responses. Variability in those
209 states leaves characteristic signatures in the statistics of sensory responses that can in principle
210 be exploited to infer states and state-specific tuning functions. We suggest that analyzing sensory
211 responses in these terms can shed new light on both their meaning and on the information they
212 carry.

213 Acknowledgements

214 We thank Alexander Ecker, Ruben Moreno-Bote, and Daniel Chicharro for feedback on the manuscript.

215 References and recommended reading

- 216 [1] Felleman, D.J. and Van Essen, D.C. (1991). Distributed hierarchical processing in the primate
217 cerebral cortex. *Cerebral Cortex* *1*, 1–47.
- 218 [2] Renart, A., Machens, C.K., and Alfonsorenart, A. (2014). Variability in neural activity and
219 behavior. *Current Opinion in Neurobiology* *25*, 211–220.
- 220 [3] Nienborg, H., Cohen, M., and Cumming, B.G. (2012). Decision-Related Activity in Sensory
221 Neurons : Correlations Among Neurons and with Behavior. *Annual Review of Neuroscience*
222 *35*, 463–483.
- 223 [4] Georgopoulos, A.P., Kettner, R.E., and Schwartz, A.B. (1988). Primate motor cortex and free
224 arm movements to visual targets in three- dimensional space. II. Coding of the direction of
225 movement by a neuronal population. *The Journal of Neuroscience* *8*, 2928–2937.
- 226 [5] Seely, J.S., Kaufman, M.T., Ryu, S.I., Shenoy, K.V., Cunningham, J.P., and Churchland, M.M.
227 (2016). Tensor Analysis Reveals Distinct Population Structure that Parallels the Different
228 Computational Roles of Areas M1 and V1. *PLOS Computational Biology* *12*, e1005164.
- 229 [6] Mitchell, J.F., Sundberg, K.A., and Reynolds, J.H. (2009). Spatial Attention Decorrelates
230 Intrinsic Activity Fluctuations in Macaque Area V4. *Neuron* *63*, 879–888.
- 231 [7] Cohen, M.R. and Maunsell, J.H.R. (2009). Attention improves performance primarily by
232 reducing interneuronal correlations. *Nature neuroscience* *12*, 1594–600.
- 233 [8] Niell, C.M. and Stryker, M.P. (2010). Modulation of visual responses by behavioral state in
234 mouse visual cortex. *Neuron* *65*, 472–479.

235 [9^{••}] Cohen, M.R. and Maunsell, J.H.R. (2010). A Neuronal Population Measure of Attention
236 Predicts Behavioral Performance on Individual Trials. *The Journal of Neuroscience* *30*, 15241–
237 15253.

238 ••Measure change in population responses as attention changes, and use this to
239 extract the magnitude of attention on a trial-by-trial basis. Inferred attention mag-
240 nitude is correlated with behavioral accuracy.

241 [10] Stănişor, L., Tógt, C.v.d., Pennartz, C.M.A., and Roelfsema, P.R. (2013). A unified selection
242 signal for attention and reward in primary visual cortex. *PNAS* *110*, 9136–9141.

243 [11^{••}] Ecker, A.S., Berens, P., Cotton, R.J., Subramaniam, M., Denfield, G.H., Cadwell, C.R.,
244 Smirnakis, S.M., Bethge, M., and Tolias, A.S. (2014). State dependence of noise correlations
245 in macaque primary visual cortex. *Neuron* *82*, 235–248.

246 ••During anesthesia, a shared source of multiplicative gain on a population changes
247 correlation structure and inflates mean correlation estimates.

248 [12] Lin, I.C., Okun, M., Carandini, M., and Harris, K.D. (2015). The Nature of Shared Cortical
249 Variability. *Neuron* *87*, 644–656.

250 [13] Schölvinck, M.L., Saleem, A.B., Benucci, A., Harris, K.D., and Carandini, M. (2015). Cortical
251 state determines global variability and correlations in visual cortex. *The Journal of Neuro-*
252 *science* *35*, 170–8.

253 [14] Nienborg, H. and Roelfsema, P.R. (2015). Belief states as a framework to explain extra-retinal
254 influences in visual cortex. *Current opinion in neurobiology* *32*, 45–52.

255 [15] Haefner, R.M., Berkes, P., and Fiser, J. (2016). Perceptual Decision-Making as Probabilistic
256 Inference by Neural Sampling. *Neuron* *90*, 649–660.

257 [16[•]] Lange, R.D. and Haefner, R.M. (2016). Inferring the brain’s internal model from sensory
258 responses in a probabilistic inference framework. *bioRxiv* .

259 •Assuming that sensory responses represent posterior beliefs, predictions are derived
260 for the structure of feedback-induced correlations. Suggest interpretation of inferred
261 internal states in terms of stimulus.

262 [17[•]] Tajima, C.I., Tajima, S., Koida, K., Komatsu, H., Aihara, K., and Suzuki, H. (2016). Popu-
263 lation code dynamics in categorical perception. *Nature Scientific Reports* *5*, 1–13.

264 •Models the effect of a categorical prior, implemented through feedback, on percep-
265 tion in a linear probabilistic population code.

266 [18[•]] Moreno-Bote, R., Beck, J.M., Kanitscheider, I., Pitkow, X., Latham, P., and Pouget, A.
267 (2014). Information-limiting correlations. *Nature Neuroscience* *17*, 1410–1417.

268 •Emphasize the point that response covariability in the $\mathbf{f}'\mathbf{f}'$ –direction limits infor-
269 mation, when induced by the feedforward pathway.

270 [19^{••}] Ecker, A.S., Denfield, G.H., Bethge, M., and Tolias, A.S. (2016). On the structure of
271 population activity under fluctuations in attentional state. *Journal of Neuroscience* *0*, 1–21.

- 272 ••A theoretical study on the effect of fluctuations in attention on response correla-
273 tions and their implication for Fisher Information.
- 274 [20] Abbott, L.F. and Dayan, P. (1999). The effect of correlated variability on the accuracy of a
275 population code. *Neural computation* *11*, 91–101.
- 276 [21] Averbek, B.B., Latham, P.E., and Pouget, A. (2006). Neural correlations, population coding
277 and computation. *Nature reviews. Neuroscience* *7*, 358–366.
- 278 [22] Bosking, W.H. and Maunsell, J.H.R. (2011). Effects of Stimulus Direction on the Correlation
279 between Behavior and Single Units in Area MT during a Motion Detection Task. *Journal of*
280 *Neuroscience* *31*, 8230–8238.
- 281 [23] Haefner, R.M. (2015). A note on choice and detect probabilities in the presence of choice bias.
282 arXiv .
- 283 [24••] Nienborg, H. and Cumming, B.G. (2009). Decision-related activity in sensory neurons reflects
284 more than a neuron’s causal effect. *Nature* *459*, 89–92.
- 285 ••Show that choice probabilities cannot be accounted for purely by a feedforward
286 readout and constant stimulus co-variability, implying choice-related activity in V2
287 contains a feedback component.
- 288 [25] Nienborg, H., Cohen, M.R., Cumming, B.G., R. Cohen, M., and Cumming, B.G. (2012).
289 Decision-related activity in sensory neurons: correlations among neurons and with behavior.
290 *Annual review of neuroscience* *35*, 463–483.
- 291 [26] Britten, K.H., Newsome, W.T., Shadlen, M.N., Celebrini, S., and Movshon, J.A. (1996). A
292 relationship between behavioral choice and the visual responses of neurons in macaque MT.
293 *Vis Neurosci* *13*, 87–100.
- 294 [27] Pitkow, X., Liu, S., Angelaki, D.E., DeAngelis, G.C., and Pouget, A. (2015). How Can Single
295 Sensory Neurons Predict Behavior? *Neuron* *87*, 411–423.
- 296 [28] Schoppik, D., Nagel, K.I., and Lisberger, S.G. (2008). Cortical Mechanisms of Smooth Eye
297 Movements Revealed by Dynamic Covariations of Neural and Behavioral Responses. *Neuron*
298 *58*, 248–260.
- 299 [29] Paninski, L., Fellows, M.R., Hatsopoulos, N.G., John, P., Wang, W., Sudre, G.P., Xu, Y.,
300 Kass, R.E., Collinger, J.L., Alan, D., et al. (2011). Spatiotemporal Tuning of Motor Cortical
301 Neurons for Hand Position and Velocity Spatiotemporal Tuning of Motor Cortical Neurons for
302 Hand Position and Velocity. *Journal of Neurophysiology* , 515–532.
- 303 [30] Haefner, R.M., Gerwinn, S., Macke, J.H., and Bethge, M. (2013). Inferring decoding strategies
304 from choice probabilities in the presence of correlated variability. *Nature neuroscience* *16*, 235–
305 42.
- 306 [31] Cohen, M.R. and Maunsell, J.H.R. (2011). Using neuronal populations to study the mecha-
307 nisms underlying spatial and feature attention. *Neuron* *70*, 1192–204.
- 308 [32•] Cunningham, J.P. and Yu, B.M. (2014). Dimensionality reduction for large-scale neural
309 recordings. *Nature Neuroscience* *17*, 1500–1509.

- 310 •Review dimensionality reduction techniques used for characterizing sensory re-
311 sponses.
- 312 [33] Putzky, P., Franzen, F., Bassetto, G., and Macke, J.H. (2014). A Bayesian model for identifying
313 hierarchically organised states in neural population activity. *Advances in Neural Information*
314 *Processing Systems 27*, 3095–3103.
- 315 [34] Kobak, D., Brendel, W., Constantinidis, C., Feierstein, C.E., Kepecs, A., Mainen, Z.F., Qi,
316 X.L., Romo, R., Uchida, N., and Machens, C.K. (2016). Demixed principal component analysis
317 of neural population data. *eLife 5*, 1–36.
- 318 [35••] Rabinowitz, N.C., Goris, R.L., Cohen, M., and Simoncelli, E.P. (2015). Attention stabilizes
319 the shared gain of V4 populations. *eLife 4*.
- 320 ••Modeled V4 neurons’ activity as a combination of known (attentional cue, stim-
321 ulus) and unknown (global gain, shared modulators) factors. Main latent factors
322 influenced responses in the same direction as the task-relevant stimulus, implying a
323 feedback source. modulators explain after known factors are accounted for.
- 324 [36] Purushothaman, G. and Bradley, D.C. (2005). Neural population code for fine perceptual
325 decisions in area MT. *Nature Neuroscience 8*, 99–106.
- 326 [37] Nienborg, H. and Cumming, B.G. (2007). Psychophysically measured task strategy for dis-
327 parity discrimination is reflected in V2 neurons. *Nature neuroscience 10*, 1608–14.
- 328 [38••] Cohen, M.R. and Newsome, W.T. (2008). Context-Dependent Changes in Functional Cir-
329 cuitry in Visual Area MT. *Neuron 60*, 162–173.
- 330 ••First to show task-dependence of response correlations in MT neurons.
- 331 [39••] Bondy, A.G. and Cumming, B.G. (2016). Feedback Dynamics Determine the Structure of
332 Spike-Count Correlation in Visual Cortex. *bioRxiv*, 1–41.
- 333 ••Fit full noise correlation matrices to macaque V1 neurons in changing task con-
334 texts. The full matrix is dominated by a task-dependent component. It is predomi-
335 nantly rank-1 and approximately proportional to $\frac{df}{dx}$ (Box 1).
- 336 [40] Markov, N.T., Ercsey-Ravasz, M., Van Essen, D.C., Knoblauch, K., Toroczkai, Z., and
337 Kennedy, H. (2013). Cortical high-density counterstream architectures. *Science 342*, 1238406.
- 338 [41] Sherman, S.M. (2016). Thalamus plays a central role in ongoing cortical functioning. *Nature*
339 *neuroscience 16*, 533–41.
- 340 [42] Briggs, F., Mangun, G.R., and Usrey, W.M. (2013). Attention enhances synaptic efficacy and
341 the signal-to-noise ratio in neural circuits. *Nature 499*, 476–80.
- 342 [43] Wimmer, K., Compte, A., Roxin, A., Peixoto, D., Renart, A., and Rocha, J.D. (2015). The
343 dynamics of sensory integration in a hierarchical network explains choice probabilities in MT.
344 *Nature Communications 6*, 1–13.
- 345 [44] Moore, T. and Zirnsak, M. (2017). Neural Mechanisms of Selective Visual Attention. *Annual*
346 *Review of Psychology 68*, 47–72.
- 347 [45] Ahissar, M. and Hochstein, S. (2004). The reverse hierarchy theory of visual perceptual
348 learning. *Trends in Cognitive Sciences 8*, 457–464.