

1 Modelling the neural code in large 2 populations of correlated neurons

3 **Sacha Sokoloski^{1*}, Amir Aschner², Ruben Coen-Cagli^{1,2}**

*For correspondence:

sacha.sokoloski@mailbox.org (SS)

4 ¹Department of Systems and Computational Biology; ²Dominick P. Purpura Department
5 of Neuroscience, Albert Einstein College of Medicine, The Bronx, New York, USA

7 **Abstract** The activity of a neural population encodes information about the stimulus that
8 caused it, and decoding population activity reveals how neural circuits process that information.
9 Correlations between neurons strongly impact both encoding and decoding, yet we still lack
10 models that simultaneously capture stimulus encoding by large populations of correlated
11 neurons and allow for accurate decoding of stimulus information, thus limiting our quantitative
12 understanding of the neural code. To address this, we propose a class of models of large-scale
13 population activity based on the theory of exponential family distributions. We apply our models
14 to macaque primary visual cortex (V1) recordings, and show they capture a wide range of
15 response statistics, facilitate accurate Bayesian decoding, and provide interpretable
16 representations of fundamental properties of the neural code. Ultimately, our framework could
17 allow researchers to quantitatively validate predictions of theories of neural coding against both
18 large-scale response recordings and cognitive performance.

20 Introduction

21 A foundational idea in sensory neuroscience is that the activity of neural populations constitutes
22 a “neural code” for representing stimuli (*Dayan and Abbott, 2005; Doya, 2007*): the activity pattern
23 of a population in response to a sensory stimulus encodes information about that stimulus, and
24 downstream neurons decode, process, and re-encode this information in their own responses.
25 Sequences of such neural populations implement the elementary functions that drive perception,
26 cognition, and behaviour (*Pitkow and Angelaki, 2017*). Therefore, by studying the encoding and de-
27 coding of population responses, researchers may investigate how information is processed along
28 neural circuits, and how this processing influences perception and behaviour (*Wei and Stocker,*
29 *2015; Panzeri et al., 2017; Kriegeskorte and Douglas, 2018*).

30 Given a true statistical model of how a neural population responds to (encodes information
31 about) stimuli, Bayes’ rule can transform the encoding model into an optimal decoder of stimu-
32 lus information (*Zemel et al., 1998; Pillow et al., 2010*). However, when validated as Bayesian de-
33 coders, existing statistical models of neural encoding are often outperformed by models trained
34 to decode stimulus-information directly, indicating that the encoding models miss key statistics of
35 the neural code (*Graf et al., 2011; Walker et al., 2020*). In particular, the correlations between neu-
36 rons’ responses to repeated presentations of a given stimulus (noise correlations), and how these
37 noise correlations are modulated by stimuli, can strongly impact coding in neural circuits (*Zohary*
38 *et al., 1994; Abbott and Dayan, 1999; Sompolinsky et al., 2001; Ecker et al., 2016; Kohn et al., 2016;*
39 *Schneidman, 2016*), especially in large populations of neurons (*Moreno-Bote et al., 2014; Montijn*
40 *et al., 2019; Bartolo et al., 2020; Kafashan et al., 2020; Rumyantsev et al., 2020*). Yet effectively
41 modelling noise correlations has proven challenging.

42 Validating theories of population coding (*Ma et al., 2006; Beck et al., 2011a; Ganguli and Simon-*

43 *celli, 2014; Makin et al., 2015; Yerxa et al., 2020*) in large neural circuits thus depends on encoding
44 models that support accurate Bayesian decoding, effectively capture noise-correlations, and effi-
45 ciently fit large-scale neural recordings. Generalized linear models (GLMs) are one class of model
46 that yield effective Bayesian decoders, and GLMs have been applied to analyzing spatio-temporal
47 features of information processing in the retina and cortex (*Pillow et al., 2008; Park et al., 2014;*
48 *Runyan et al., 2017*). Nevertheless, neural correlations are often the result of low-dimensional,
49 shared variability (*Arieli et al., 1996; Ecker et al., 2014; Goris et al., 2014; Rabinowitz et al., 2015;*
50 *Okun et al., 2015; Semedo et al., 2019*), and it is unknown whether extensions of the GLM ap-
51 proach to capture shared-variability (*Archer et al., 2014; Zhao and Park, 2017*) can support accu-
52 rate Bayesian decoding. Similarly, methods based on factor analysis (*Yu et al., 2009; Ecker et al.,*
53 *2014; Semedo et al., 2019*) have proven highly effective at modelling neural correlations in large-
54 scale recordings, but it also unknown if they can support Bayesian decoding. Finally, a model class
55 related to GLMs is pairwise-maximum entropy models (*Schneidman et al., 2006; Lyamzin et al.,*
56 *2010; Granot-Atedgi et al., 2013; Meshulam et al., 2017*), which have been used to investigate se-
57 mantic clustering of responses in the retinal code (*Ganmor et al., 2015*); yet these models have
58 so-far been limited to population sizes of tens of neurons.

59 Towards modelling responses and accurate Bayesian decoding in large populations of corre-
60 lated neurons, we have developed a class of spike-count encoding model based on conditional fi-
61 nite mixtures of multivariate Poisson distributions, which we refer to as CPMs (Conditional Poisson
62 Mixtures). Within neuroscience, Poisson mixtures are widely applied to modelling the spike-count
63 distributions of individual neurons (*Wiener and Richmond, 2003; Shidara et al., 2005; Goris et al.,*
64 *2014; Taouali et al., 2015*). Outside of neuroscience, mixtures of multivariate Poisson distributions
65 are an established model of multivariate count distributions that effectively capture correlations
66 in count data (*Karlis and Meligkotsidou, 2007; Inouye et al., 2017*).

67 Building on the theory of exponential family distributions (*Wainwright and Jordan, 2008; Macke*
68 *et al., 2011b*), our model extends previous mixture models of multivariate count data in two ways.
69 Firstly, we develop a tractable extension of Poisson mixtures that captures both over- and under-
70 dispersed response variability (i.e. where the response variance is larger or smaller than the mean,
71 respectively) based on Conway-Maxwell Poisson distributions (*Shmueli et al., 2005; Stevenson,*
72 *2016*). Secondly, we introduce an explicit dependence of the model on a stimulus variable, which
73 allows the model to accurately capture changes in response statistics (including noise correlations)
74 across stimuli. Importantly, the resulting encoding model affords closed-form expressions for both
75 its Fisher information and probability density function and thereby a rigorous quantification of
76 the coding properties of a modelled neural population (*Dayan and Abbott, 2005*). Moreover, the
77 model learns low-dimensional representations of stimulus-driven neural activity, and we show how
78 it provides a parsimonious description of a fundamental property of population codes known as
79 information-limiting correlations (*Moreno-Bote et al., 2014; Montijn et al., 2019; Bartolo et al.,*
80 *2020; Kafashan et al., 2020; Rumyantsev et al., 2020*).

81 We apply the CPM framework to both synthetic data and recordings from macaque primary
82 visual cortex (V1), and demonstrate that it effectively models responses of populations of hundreds
83 of neurons, captures noise correlations, and supports accurate Bayesian decoding. Ultimately, our
84 model of neural encoding and decoding can be used to quantify coding properties of a neural
85 circuit, such as their efficiency, linearity, or information capacity.

86 Results

87 A critical part of our theoretical approach is based on expressing models of interest in exponen-
88 tial family form. An exponential family distribution $p(n)$ over some data n (in our case, neural re-
89 sponses) is defined by the proportionality relation $p(n) \propto e^{\theta \cdot s(n)} b(n)$, where θ are the so-called natural
90 parameters, $s(n)$ is a vector-valued function of the data called the sufficient statistic, and $b(n)$ is a
91 scalar-valued function called the base measure (*Wainwright and Jordan, 2008*). The exponential
92 family form allows us to modify and extend existing models in a simple and flexible manner, and

93 to gain analytical insight into the coding properties of our models. We demonstrate our approach
94 with applications to both synthetic data generated by example CPMs, and data recorded in V1 of
95 anaesthetized and awake macaques viewing drifting grating stimuli at different orientations (for
96 details see Materials and methods).

97 **Extended Poisson mixture models capture spike-count variability and covariability**

98 Our first goal is to define a class of models of neural population activity, that model neural activity
99 directly as spike-counts, and that accurately capture single-neuron variability and pairwise covari-
100 ability. We base our models on Poisson distributions, as they are widely-applied to modelling the
101 trial-to-trial distribution of the number of spikes generated by a neuron (*Dayan and Abbott, 2005*;
102 *Macke et al., 2011a*). We will also generalize our Poisson models with Conway-Maxwell (CoM) Pois-
103 son distributions, because they can capture the broad range of Fano factors (FF; the variance di-
104 vided by the mean) observed in cortex, in contrast with Poisson distributions for which the FF is
105 always 1 (*Sur et al., 2015*; *Stevenson, 2016*; *Chaniialidis et al., 2018*).

106 Mixtures of Poisson distributions are also used to capture complex spike-count distributions
107 in cortex, and allow for over-dispersion (FF>1) (*Shidara et al., 2005*; *Goris et al., 2014*; *Taouali*
108 *et al., 2015*) (Figure 1A). In our case we consider multivariate Poisson mixtures, as they capture
109 covariability in count data as well (see *Karlis and Meligkotsidou (2007)* for the general definition).
110 To construct a multivariate Poisson mixture we begin with a product of independent Poisson dis-
111 tributions, one per neuron. We then mix a finite number of such independent Poisson models,
112 to arrive at a multivariate spike-count, finite mixture model (see Materials and methods). Impor-
113 tantly, although each mixture component is a product of independent Poisson distributions, ran-
114 domly switching between components induces correlations between the neurons (Figure 1B,C). In
115 fact, multivariate Poisson mixtures may model arbitrary pairwise covariability (see Materials and
116 methods, Equation 6). Nevertheless, they are limited because the variance of individual neurons
117 cannot be smaller than the mean, and are thus always over-dispersed (Equation 5, Materials and
118 methods).

119 To address this limitation, we show how to express multivariate Poisson mixtures in an expo-
120 nential family form, and then generalize the model with CoM-Poisson distributions. We first note
121 that a multivariate Poisson mixture with d_K components may be expressed as a latent variable
122 model over spike-count vectors \mathbf{n} and latent component-indices k , where $1 \leq k \leq d_K$. In this formu-
123 lation we denote the k th component distribution by $p(\mathbf{n} | k)$, and the probability of realizing (switch-
124 ing to) the k th component by $p(k)$. The mixture model over spike-counts \mathbf{n} is then expressed as
125 the marginal distribution $p(\mathbf{n}) = \sum_{k=1}^{d_K} p(\mathbf{n} | k)p(k) = \sum_{k=1}^{d_K} p(\mathbf{n}, k)$, of the joint distribution $p(\mathbf{n}, k)$. Un-
126 der mild regularity assumptions (see Materials and methods), we may reparameterize this joint
127 distribution in exponential family form as

$$p(\mathbf{n}, k) \propto \frac{e^{\theta_N \cdot \mathbf{n} + \theta_K \cdot \delta(k) + \Theta_{NK} \cdot \delta(k)}}{\prod_{i=1}^{d_N} n_i!}, \quad (1)$$

128 where θ_N , θ_K , and Θ_{NK} are the natural parameters of $p(\mathbf{n}, k)$, and $\delta(k) = (\delta_2(k), \dots, \delta_{d_K}(k))$ is the
129 Kronecker delta vector defined by $\delta_j(k) = 1$ if $j = k$, and 0 otherwise.

130 The exponential family form of a multivariate Poisson mixture represents the first component
131 distribution (i.e. $p(\mathbf{n} | k)$ with index $k = 1$) as a baseline distribution, and the other components
132 (where $k > 1$) as modulations of the baseline distribution, and this representation helps us extend
133 multivariate Poisson mixtures. In particular, the first component distribution has natural (base-
134 line) parameters θ_N , and for $k > 1$, the natural parameters of $p(\mathbf{n} | k)$ are the sum of the baseline
135 parameters θ_N and one row from the matrix of parameters Θ_{NK} (Equation 12, Materials and meth-
136 ods). Because the dimension of θ_N is much smaller than the total number of parameters in a given
137 mixture, the baseline parameters provide a relatively low-dimensional means of affecting all the
138 component distributions of the given mixture, as well as the index probabilities (Figure 1D; see
139 Materials and methods, Equation 11 for how $p(k)$ depends on θ_N).

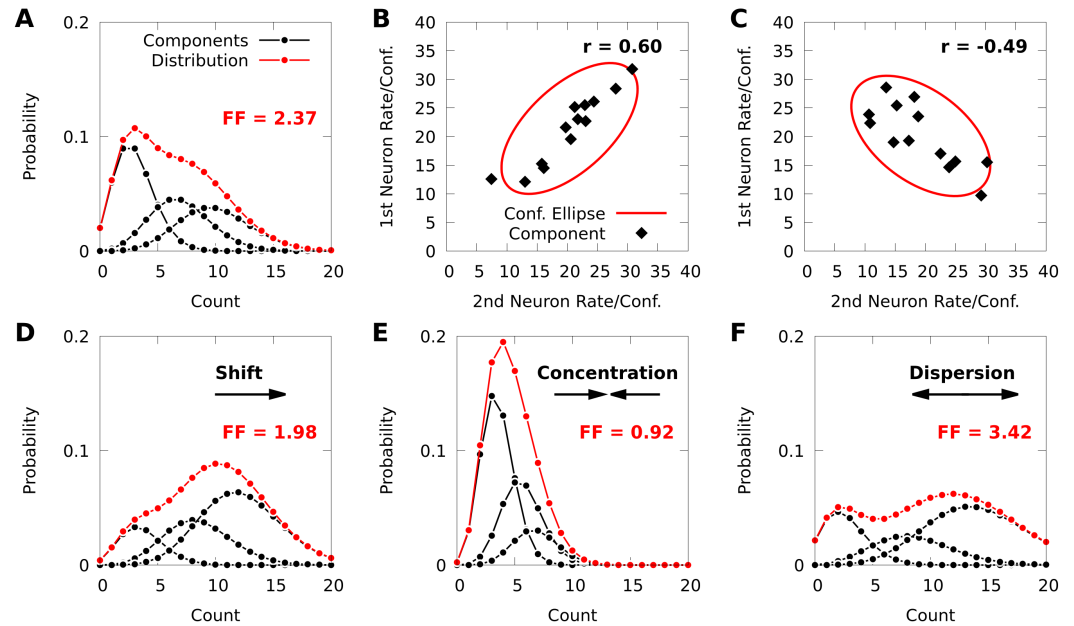


Figure 1. *Poisson mixtures and Conway-Maxwell extensions* **A:** A Poisson mixture distribution (red), defined as the weighted sum of three component Poisson distributions (black; scaled by their weights). FF denotes the Fano Factor (variance over mean) of the mixture. **B,C:** The average spike-count (rate) of the first and second neurons for each of 13 components (black dots) of a bivariate Poisson mixture model, and 68% confidence ellipses for the spike-count covariance of the mixture (red lines; see Equations 5 and 6). The spike-count correlation of each mixture is denoted by r . **D:** Same model as **A**, except we shift the distribution by increasing the baseline rate of the components. **E,F:** Same model as **A**, except we use an additional baseline parameter based on Conway-Maxwell Poisson distributions to concentrate (**E**) or disperse (**F**) the mixture distribution and its components.

140 We now extend Relation 1 with CoM-Poisson theory, and propose the latent variable exponential
141 family

$$p(\mathbf{n}, k) \propto e^{\theta_N \cdot \mathbf{n} + \theta_N^* \cdot \mathbf{If}(\mathbf{n}) + \theta_K \cdot \delta(k) + \mathbf{n} \cdot \Theta_{NK} \cdot \delta(k)}, \quad (2)$$

142 where $\mathbf{If}(\mathbf{n}) = (\log(n_1!), \dots, \log(n_{d_N}!))$ is the vector of log-factorials of the individual spike-counts,
143 and θ_N^* are a set of natural parameters based on CoM-Poisson distributions (see Materials and
144 methods). The exponential family form continues to represent the mixture in terms of a baseline
145 distribution, in this case $p(\mathbf{n} | k)$ is a product of independent CoM-Poisson distributions, with base-
146 line parameters θ_N and CoM-based parameters θ_N^* . However, whereas the rows of Θ_{NK} modulate
147 θ_N depending on the component index k , the parameters θ_N^* are not modulated, and remain the
148 same for each component distribution (Equation 15, Materials and methods, and see Equation 14
149 for index-probability formula). For the rest of this paper we refer to models described by Relation 1
150 as vanilla mixtures, and models described by Relation 2 as CoM-based mixtures.

151 Due to the addition of the CoM-based parameters, a CoM-based mixture can model under-
152 dispersed ($FF < 1$) neural activity (Equation 16, Materials and methods). In Figures 1D-F we demon-
153 strate how changing the parameters of the CoM-based mixture can concentrate or disperse both
154 the mixture distribution and its components.

155 To validate our mixture models, we tested if they capture variability and covariability of V1 pop-
156 ulation responses to repeated presentations of a grating stimulus with fixed orientation ($d_N = 43$
157 neurons and $d_T = 355$ repetitions in one awake macaque; $d_N = 70$ and $d_T = 1,200$ in one anaes-
158 thetized macaque). We optimized model parameters as described in Materials and methods. The
159 CoM-Poisson mixture accurately captured single-neuron variability (Figure 2A-B, red symbols), in-
160 cluding both cases of over-dispersion and under-dispersion. In contrast, the simpler multivariate
161 Poisson mixture (Figure 2A-B, blue symbols) could not accommodate under-dispersion, and also

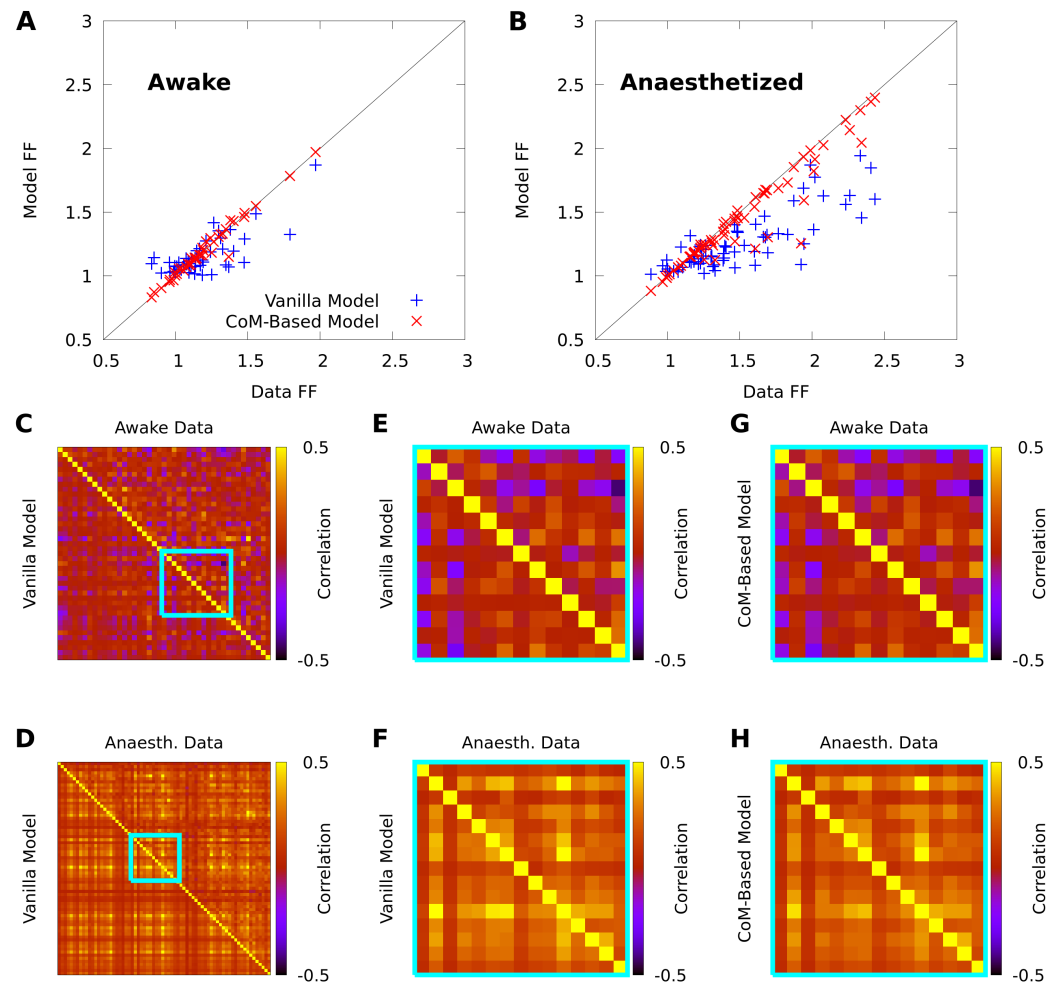


Figure 2. Capturing neural variability in V1 responses to a single stimulus with CPMs. We qualitatively compare vanilla Poisson mixtures (Relation 1) and CoM-based mixtures (Relation 2), on awake and anaesthetized V1 responses to stimulus orientation $x = 20^\circ$; both mixtures are defined with $d_K = 4$ components for awake data, and $d_K = 8$ components for anaesthetized data (see Materials and methods for training algorithms). **A,B:** Empirical Fano factors of the awake (**A**) and anaesthetized data (**B**), compared to vanilla (blue) and CoM-based mixtures (red). **C,D:** Empirical correlation matrix (upper right) of awake (**C**) and anaesthetized data (**D**), compared to the correlation matrix of the corresponding vanilla mixtures (lower left). **E,F:** Correlations highlighted in **C** and **D**, respectively. **G,H:** Correlations highlighted in **C** and **D**, except model correlations are from CoM-based mixtures.

162 had a limited ability to model over-dispersion due to the coupling between the mean and variance
 163 (Equation 5). On the other hand, we found that both mixture models were flexible enough to qual-
 164 itatively capture pairwise noise correlations, both in awake and anaesthetized animals (Fig. 2C-H)
 165 (in later sections we quantitatively compare the model performance).

166 Extended Poisson mixture models capture stimulus-dependent response statistics

167 So far we have introduced the exponential family theory of vanilla and CoM-based Poisson mix-
 168 tures, and shown how they capture response variability and covariability for a fixed stimulus. To
 169 allow us to study stimulus encoding and decoding, we further extend our mixtures by inducing
 170 a dependency of the model parameters on a stimulus. When there are a finite number of stim-
 171 ulus conditions and sufficient data, we may define a stimulus-dependent model with a lookup
 172 table, and fit it by fitting a distinct model at each stimulus condition. However, this is inefficient

173 when the amount of data at each stimulus-condition is limited and the stimulus-dependent statis-
174 tics have structure that is shared across conditions. A notable feature of the exponential family
175 parameterizations in Relations 1 and 2 is that the baseline parameters influence both the index
176 probabilities and all the component distributions of the model. This suggests that by restricting
177 stimulus-dependence to the baseline parameters, we might model rich stimulus-dependent re-
178 sponse structure, while bounding the complexity of the model.

179 In general we refer to any finite mixture of independent Poisson distributions with stimulus-
180 dependent parameters as a conditional Poisson mixture (CPM), and depending on whether the
181 CPM is based on Relations 1 or 2, we refer to it as a vanilla or CoM-based CPM, respectively. Al-
182 though there are many ways we might induce stimulus-dependence, in this paper we consider
183 two forms of CPM: (i) a maximal CPM, which we implement as a lookup table, such that all the
184 parameters in Relation 1 or 2 depend on the stimulus, and (ii) a minimal CPM, for which we restrict
185 stimulus-dependence to the baseline parameters θ_N , resulting in the CoM-based CPM

$$p(\mathbf{n}, k | x) \propto e^{\theta_N(x) \cdot \mathbf{n} + \theta_N^* \cdot \mathbf{1}(\mathbf{n}) + \theta_K \cdot \delta(k) + \mathbf{n} \cdot \Theta_{NK} \cdot \delta(k)}, \quad (3)$$

186 where x is the stimulus, and $\theta_N(x)$ are the stimulus-dependent baseline parameters (we may re-
187 cover a minimal, vanilla CPM by setting $\theta_N^* = -\mathbf{1}$). The tuning curves of the CPM neurons are the
188 average spike-counts (firing rates) of each n_i as a function of the stimulus x , and we refer to $\theta_N(x)$
189 as the baseline tuning curve parameters, as they define how the firing rates of the baseline CPM
190 distribution (i.e. $p(\mathbf{n} | x, k)$ when $k = 1$) depend on x . For $k > 1$, the modulated CPM $p(\mathbf{n} | x, k)$ is
191 then a scaled, or “gain-modulated” version of the baseline CPM (see Equations 12 and 15 and the
192 accompanying discussions).

193 Towards understanding the expressive power of CPMs, we study a minimal, CoM-based CPM
194 with $d_N = 20$ neurons, $d_K = 5$ mixture components, and randomly chosen parameters (see Ma-
195 terials and methods). Moreover, we assume that the stimulus is periodic (e.g. the orientation of
196 a grating), and that the baseline tuning curves have a von Mises shape which is a widely applied
197 model of neural tuning to periodic stimuli (*Herz et al., 2017*). We may achieve such a baseline
198 shape by defining the baseline tuning curve parameters as $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \mathbf{vm}(x)$, where θ_N^0
199 and Θ_{NX} are the tuning curve parameters, and $\mathbf{vm}(x) = (\cos 2x, \sin 2x)$. Figure 3A shows that the
200 tuning curves of the CPM neurons are approximately bell-shaped, yet many also exhibit significant
201 deviations.

202 We also study if CPMs can be effectively fit to datasets comparable to those obtained in typical
203 neurophysiology experiments. We generated 200 responses from the CoM-based CPM described
204 above — the ground truth CPM — to each of 10 orientations spread evenly over the half-circle, for a
205 total of 2,000 stimulus-response sample points. We then used this data to fit a CPM with the same
206 number of components. Towards this aim, we derived an approximate expectation-maximization
207 algorithm (EM, a standard choice for training finite mixture models (*McLachlan et al., 2019*)) to
208 optimize model parameters, that also accounts for the stimulus-dependence (see Materials and
209 methods). Figure 3B shows that the tuning curves of the learned CPM are nearly indistinguishable
210 from those of the ground truth CPM (Figure 3B).

211 To reveal the orientation-dependent latent structure of the model, in Figure 3C we plot the
212 index probability $p(k | x)$ for every k as a function of the orientation x . In Figure 3D we show
213 that the orientation-dependent index probabilities of the learned CPM qualitatively match the true
214 index probabilities in Figure 3C. We also note that although the learned CPM does not correctly
215 identify the indices themselves, this has no effect on the performance of the CPM.

216 The orientation-dependent index-probabilities provide a high-level picture of how the complex-
217 ity and structure of model correlations varies with the orientation. The vertical dashed lines in
218 Figures 3C-D denote two orientations that yield substantially different index probabilities $p(k | x)$.
219 When a large number of index-probabilities are non-zero, the correlation-matrices of the CoM-
220 based CPM can exhibit complex correlations with both negative and positive values (Figure 3E).
221 However, when one index dominates, the correlation structure largely disappears (Figure 3F). In

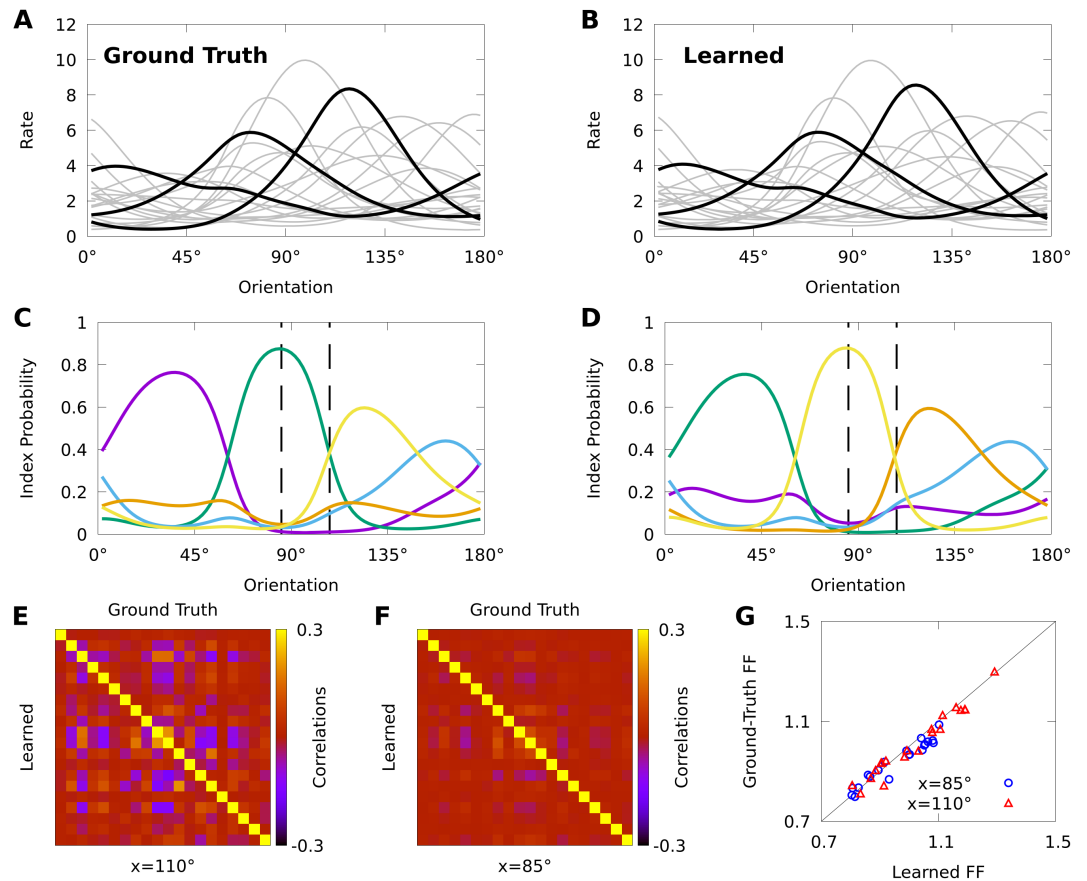


Figure 3. Recovering a ground truth conditional Poisson mixture (CPM). We compare a ground truth, CoM-based CPM with 20 neurons, 5 mixture components, von Mises baseline tuning, and randomized parameters to a learned CPM fit to 2,000 samples from the ground truth CPM. **A-B:** Tuning curves of the ground-truth CPM (**A**) and learned CPM (**B**). Three tuning curves are highlighted for effect. **C-D:** The orientation-dependent index probabilities of the ground truth CPM (**C**) and learned CPM (**D**), where colour indicates component index. Dashed lines indicate example stimulus-orientations used in Figures 3C-D. **E-F:** The correlation matrix of the ground truth CPM (upper right), compared to the correlation matrix of the learned CPM (lower left) at stimulus orientations $x = 85^\circ$ (**E**) and $x = 110^\circ$ (**F**). **G:** The FFs of the ground-truth CPM compared to the learned CPM at orientations $x = 85^\circ$ (blue circles) and $x = 110^\circ$ (red triangles).

222 Figure 3G we show that the FFs also depend on stimulus orientation. Lastly, we find that both the
 223 FF and the correlation-matrices of the learned CPM are nearly indistinguishable from the ground-
 224 truth CPM (Figure 3E-G).

225 In summary, our analyses show that CPMs can generate complex, stimulus-dependent response
 226 statistics, and that the learned CPM accurately recovers both the statistics and the latent structure
 227 of the neural responses from realistic amounts of data.

228 CPMs effectively model neural responses in macaque V1

229 A variety of models may be defined within the CPM framework illustrated by Relations 1, 2, and 3.
 230 Towards understanding how effectively CPMs can model real data, we compare different variants
 231 by their cross-validated log-likelihood. We consider both vanilla and CoM-based variants of each of
 232 the following conditional mixtures: (i) maximal CPMs where we learn a distinct mixture for each of
 233 d_x stimulus conditions, (ii) minimal CPMs with von Mises baseline tuning curves, and (iii) minimal
 234 CPMs with *discrete* baseline tuning curves given by $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \delta(x)$, where δ is the Kronecker
 235 delta vector with $d_x - 1$ elements, and x is the index of the stimulus. In contrast with the von Mises
 236 CPM, the discrete CPM makes no assumptions about the form of baseline tuning.

Encoding Performance

	Awake			Anaesthetized		
	Inf. Gain	d_K	Num. Params.	Inf. Gain	d_K	Num. Params.
Maximal Vanilla	2.30 ± 0.32	5	2,689	8.77 ± 0.71	8	5,103
Maximal CoM	2.44 ± 0.35	5	3,044	9.42 ± 0.70	7	4,464
VM Vanilla	2.01 ± 0.26	45	2,065	8.97 ± 0.70	40	2,979
VM CoM	2.10 ± 0.25	40	1,888	9.38 ± 0.69	35	2,694
Disrete Vanilla	2.25 ± 0.28	40	2,103	9.17 ± 0.70	35	3,044
Disrete CoM	2.35 ± 0.29	30	1,708	9.53 ± 0.68	30	2,689

Table 1. The encoding performance of CPMs on neural responses in macaque V1. We apply 10-fold cross-validation to estimate the mean and standard error of the information gain on held-out data, from either awake or anaesthetized macaque V1. We compare maximal CPMs (Maximal), minimal CPMs with von Mises baseline tuning (VM), and minimal CPMs with discrete baseline tuning (Discrete), and for each case we consider either Vanilla or CoM-based variants. For each variant, we indicate the number of CPM components d_K and the corresponding number of model parameters required to achieve peak information gain (cross-validated). For reference, the independent Poisson models use 129 and 210 parameters for the awake and anaesthetized data, respectively.

237 To provide an interpretable measure of the relative performance of each CPM variant, we mea-
 238 sured the difference between the estimated log-likelihood of the given CPM and the log-likelihood
 239 of a von Mises-tuned, independent Poisson model, which is a standard model of uncorrelated neu-
 240 ral responses to oriented stimuli (Herz et al., 2017). We refer to this quantity as the information
 241 gain.

242 Table 1 shows that the CPM variants considered achieve comparable performance, and per-
 243 form substantially better than the independent Poisson lower bound on both the awake and anaes-
 244 thetized data. Figure 4 shows that a performance peak emerges smoothly as the model complexity
 245 (number of parameters) is increased. In all cases, the CoM-based models outperform their vanilla
 246 counterparts, and typically with fewer parameters. The CoM-based discrete CPMs achieve high
 247 performance on both datasets. In contrast, von Mises CPMs perform well on the anaesthetized
 248 data but more poorly on the awake data, and maximal CPMs exhibit the opposite trend. Never-
 249 theless, von Mises CPMs solve a more difficult statistical problem as they also interpolate between
 250 stimulus conditions, and so may still prove relevant even where performance is limited. On the
 251 other hand, even though maximal CPMs achieve high performance, they simply do so by replicat-
 252 ing the high performance of stimulus-independent mixtures (Figure 2) at each stimulus condition,
 253 requiring significantly more parameters than minimal CPMs.

254 CPMs facilitate accurate and efficient Bayesian decoding of neural responses

255 To demonstrate that CPMs model the neural code, we must show that CPMs not only capture the
 256 features of neural responses, but that these features also encode stimulus-information. Given an
 257 encoding model $p(\mathbf{n} | x)$ and a response from the model \mathbf{n} , we may optimally decode the informa-
 258 tion in the response about the stimulus x by applying Bayes' rule $p(x | \mathbf{n}) \propto p(\mathbf{n} | x)p(x)$, where $p(x | \mathbf{n})$
 259 is the posterior distribution (the decoded information), and $p(x)$ represents our prior assumptions
 260 about the stimulus (Zemel et al., 1998). When we do not know the true encoding model, and rather
 261 fit a statistical model to stimulus-response data, using the statistical model for Bayesian decoding
 262 and analyzing its performance can tell us how well it captures the features of the neural code.

263 We analyze the performance of Bayesian decoders based on CPMs by quantifying their decod-
 264 ing performance, and comparing the results to other common approaches to decoding. We quan-
 265 tify decoding performance by evaluating the average of the cross-validated log-posterior probabili-
 266 ty $\log p(x | \mathbf{n})$ of the true stimulus value x , for both our awake and anaesthetized V1 datasets. With
 267 regards to training the CPMs, we analyze the decoding performance of CPMs that achieved the

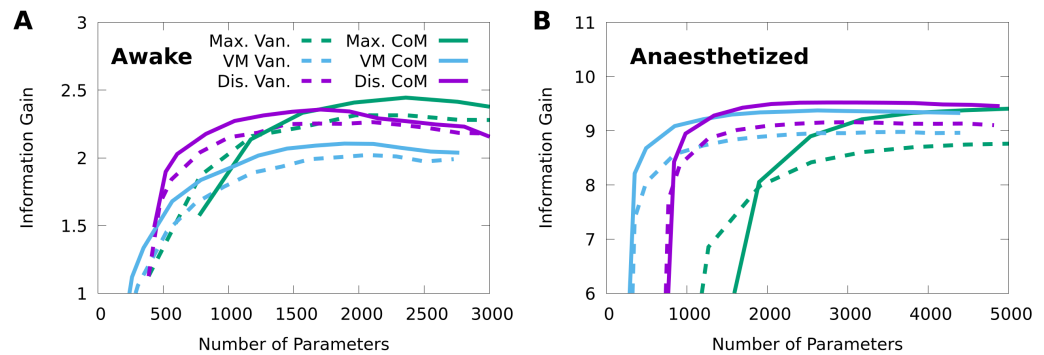


Figure 4. Finding the optimal number of parameters for CPMs to model neural responses in macaque V1. 10-fold cross-validation of the information gain given awake V1 data (A) and anaesthetized V1 data (B), as a function of the number of model parameters, for multiple forms of CPM: maximal CPMs (green); minimal CPMs with von Mises baseline tuning (blue); minimal CPMs with discrete baseline tuning (purple); and for each case we consider either vanilla (dashed lines) or CoM-based (solid lines). Standard errors of the information gain are not depicted to avoid visual clutter, however they are approximately independent of the number of model parameters, and match the values indicated in Table 1.

268 best *encoding* performance based as indicated in Table 1 and depicted Figure 4, instead of apply-
269 ing distinct procedures for selecting CPMs based on decoding performance. This is because our
270 goal is to understand how well the response features captured by CPMs reflect the neural code,
271 rather than strictly maximizing decoding performance.

272 In our comparisons we focus on minimal, discrete CPMs as overall they achieved high per-
273 formance on both datasets (Figure 4). To characterize the importance of neural correlations to
274 Bayesian decoding, we compare our CPMs to the decoding performance of independent Poisson
275 models with discrete tuning (IP). To characterize the optimality of our Bayesian decoders, we also
276 evaluate the performance of linear multiclass decoders (Linear), as well nonlinear multiclass de-
277 coders defined as artificial neural networks (ANNs) with two hidden layers and a cross-validated
278 number of hidden units (for details on the training and model selection procedure, see *Materials*
279 *and methods*).

280 Table 2 shows that on the awake data, the performance of the CPMs is statistically indistinguish-
281 able from the ANN, and the CPMs and the ANN significantly exceed the performance of both the
282 Linear and IP models. On the anaesthetized data, the minimal CPM approaches the performance
283 of the ANN, and the minimal CPMs and ANN models again exceed the performance of the IP and
284 Linear models. Yet in this case the Linear model is much more competitive, whereas the IP model
285 performs very poorly, possibly because of the larger magnitude of noise correlations in this data. In
286 both cases the ANN requires two orders of magnitude more parameters than the CPMs to achieve
287 its performance gains. In addition, the CoM-based CPM achieves marginally better performance
288 with fewer parameters than the vanilla CPM, indicating that although modelling individual variabil-
289 ity is not essential for effective Bayesian decoding, doing so still results in a more parsimonious
290 model of the neural code.

291 We also consider widely used alternative measures of decoding performance, namely the Fisher
292 information (FI), which is an upper bound on the average precision (inverse variance) of the pos-
293 terior (*Brunel and Nadal, 1998*), as well as the linear Fisher information (LFI), which is a linear ap-
294 proximation of the FI (*Seriès et al., 2004*) corresponding to the accuracy of the optimal, unbiased
295 linear decoder of the stimulus (*Kanitscheider et al., 2015a*). The FI is especially helpful when the
296 posterior cannot be evaluated directly (such as when it is continuous), and is widely adopted in the-
297 oretical (*Abbott and Dayan, 1999; Ecker et al., 2014; Moreno-Bote et al., 2014; Kohn et al., 2016*)
298 and experimental (*Ecker et al., 2011; Rumyantsev et al., 2020*) studies of neural coding. As with
299 other models based on exponential family theory (*Ma et al., 2006; Beck et al., 2011b; Ecker et al.,*
300 *2016*), the FI of a minimal CPM may be expressed in closed-form, and is equal to its LFI (see Ma-

Decoding Performance

	Awake		Anaesthetized	
	Average Log-Post.	Num. Params.	Average Log-Post.	Num. Params.
CoM CPM	-0.206 ± 0.043	1,663	-0.441 ± 0.023	2,689
Vanilla CPM	-0.207 ± 0.039	2,103	-0.448 ± 0.026	3,044
Ind. Poisson	-0.272 ± 0.067	387	-0.967 ± 0.071	630
Linear	-0.256 ± 0.053	352	-0.457 ± 0.019	568
Artificial NN	-0.200 ± 0.032	527,108	-0.426 ± 0.015	408,008

Table 2. The decoding performance of CPMs on neural responses in macaque V1. We apply 10-fold cross-validation to estimate the mean and standard error of the average log-posteriors $\log p(x | \mathbf{n})$ on held-out data, from either awake or anaesthetized macaque V1. We compare discrete, minimal, CoM-based CPM (CoM CPM) and vanilla CPM (Vanilla CPM); an independent Poisson model with discrete tuning (IP); a multiclass linear decoder (Linear); and a multiclass nonlinear decoder defined as an artificial neural network with two hidden layers (ANN). The number of CPM components d_K was chosen to achieve on peak information gain in Figure 4. The number of ANN hidden units was chosen based on peak cross-validation performance. In all cases we also indicate the number of model parameters required to achieve the indicated performance.

301 terials and methods), and therefore minimal CPMs can be used to study FI analytically and obtain
302 model-based estimates of FI from data.

303 We generated 40 populations of $d_N = 20$ model neurons from the vanilla, minimal, von Mises
304 CPM, with parameters corresponding to the best-fit parameters of 40 random subsets of neurons
305 from our V1 datasets. For each population, we generated 50 responses at each of 10 evenly spaced
306 orientations, for a total of $d_T = 500$ responses per population. We then fit a CPM to each set of
307 500 responses, and compared the FI of the fit CPM to the ground-truth FI at 50 evenly spaced
308 orientations. Pooled over all populations and orientations, the relative error of the estimated FI
309 was $-12.8\% \pm 18.6\%$ on the awake data and $-9.1\% \pm 22.4\%$ on the anaesthetized data.

310 The aforementioned measures allow us to assess decoding performance when we do not know
311 the full posterior, however the full posterior is an essential part of probabilistic neural codes (Pouget
312 *et al.*, 2016; Drugowitsch *et al.*, 2019). To test whether CPMs can in principle recover full posteri-
313 ors, we consider a ground truth CPM defined as discrete, CoM-based, minimal CPM with $d_N = 200$
314 neurons, $d_S = 20$ stimulus-conditions, $d_K = 30$ components, and randomized parameters, and we
315 fit a discrete, CoM-based, minimal CPM with $d_K = 40$ components (chosen with cross-validation) to
316 $d_T = 10,000$ responses from the ground-truth CPM (see Materials and methods). We then compute
317 the average KL-divergence (a fundamental measure of the similarity of two distributions, see Cover
318 and Thomas (2006); Amari and Nagaoka (2007)) of the learned posteriors from the ground-truth
319 posterior over all the $d_T = 10,000$ responses, and find that the average posterior divergence is
320 0.047 ± 0.007 bits, indicating that on average the learned and ground-truth posteriors are extremely
321 close.

322 To summarize, CPMs support accurate Bayesian decoding in awake and anaesthetized macaque
323 V1 recordings, and are competitive with nonlinear decoders with two orders of magnitude more
324 parameters. Moreover, CPMs afford closed-form expressions of FI and can interpolate good esti-
325 mates of FI from modest amounts of data, and thereby support analyses of neural data based on
326 this widely applied theoretical tool. Finally, a CPM fit to the responses of a ground-truth CPM can
327 almost perfectly recover the ground-truth posterior distributions.

328 Minimal CPMs provide an interpretable latent representation of a fundamental 329 feature of the neural code

330 Having shown that CPMs can be used to accurately decode stimuli, we next aim to demonstrate
331 that the latent structure of CPMs offers an interpretable representation of a central phenomenon
332 in neural coding known as information-limiting correlations, which are neural correlations that

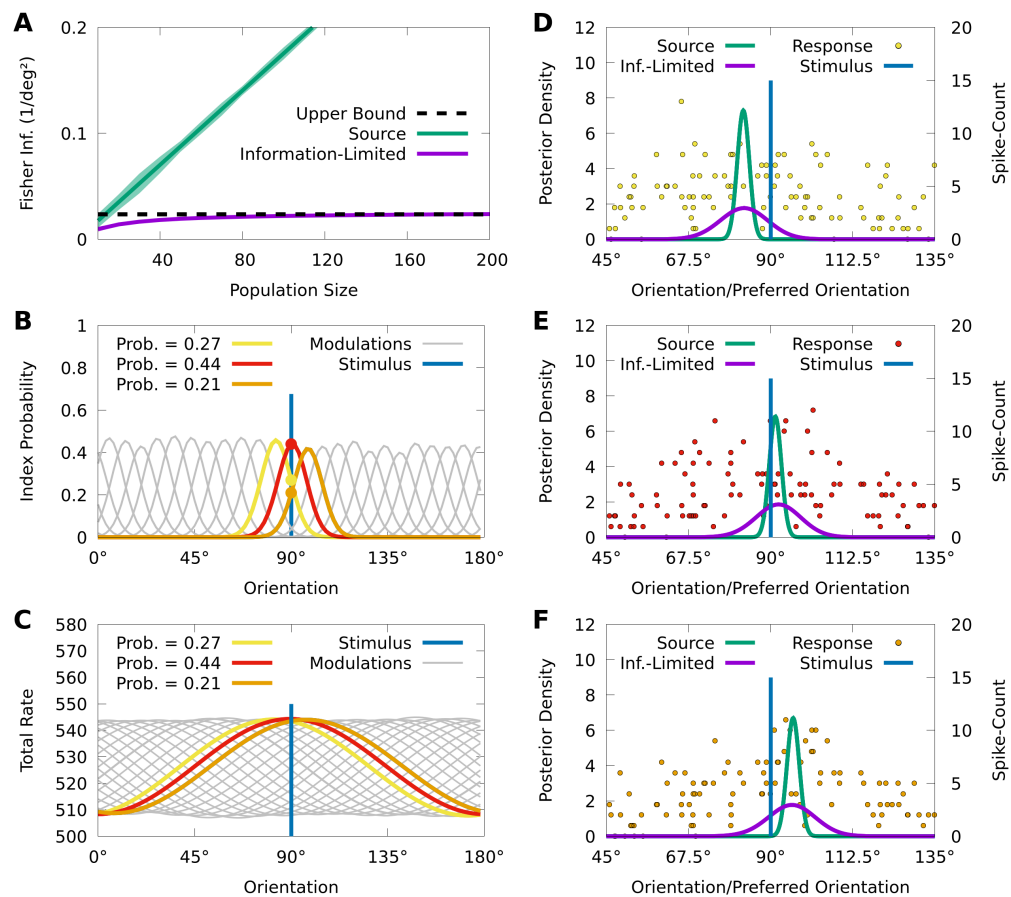


Figure 5. Fisher information and information-limiting correlations in CPMS. We consider a von Mises-tuned, independent Poisson source model (green) with $d_K = 200$ neurons, and an information-limited, CoM-based CPM (purple) with $d_K = 25$ components, fit to 10,000 responses of the source-model to stimuli obscured by von Mises noise. In **B-F** we consider a stimulus-orientation $x = 90^\circ$ (blue line). **A**: The average (lines) and standard deviation (filled area) of the FI over orientations, for the source (green) and information-limited (purple) models, as a function of random subpopulations, starting with ten neurons, and gradually reintroducing missing neurons. Dashed black line indicates the theoretical upper bound. **B**: The index-probability curves (lines) of the CPM for indices $k > 1$ and the intersection (red, yellow, and orange circles) of the stimulus with three curves (orange, yellow, and orange lines). **C**: The sum of the firing rates of the modulated CPM for all indices $k > 1$ (lines) as a function of orientation, with three modulated CPMS highlighted (red, yellow, and orange lines) corresponding to the highlighted indices in **B**. **D-F**: Three responses from the yellow (**D**; yellow points), red (**E**; red points), and orange modulated CPMS (**F**; orange points) indicated in **C**. For each response we plot the posterior based on the source model (green line) and the information-limited model (purple line).

333 fundamentally limit stimulus-information in neural circuits (*Moreno-Bote et al., 2014; Montijn et al.,*
 334 *2019; Bartolo et al., 2020; Kafashan et al., 2020; Rumyantsev et al., 2020*). To illustrate this, we
 335 generate population responses with limited information, and then fit a CPM to these responses
 336 and study the learned latent representation. In particular, we consider a source population of 200
 337 independent Poisson neurons $p(\mathbf{n} | s)$ with homogeneous, von Mises tuning curves responding to a
 338 noisy stimulus-orientation s , where the noise $p(s | x)$ follows a von Mises distribution centred at the
 339 true stimulus-orientation x (see *Materials and methods*). In Figure 5A we show that, as expected,
 340 the average FI in the source population about the noisy orientation s grows linearly with the size of
 341 randomized subpopulations, whereas the FI about the true orientation x is theoretically bounded
 342 by the precision (inverse variance) of the sensory noise.

343 Even though the neurons in the source model are uncorrelated, sensory noise ensures that

344 the information-limited encoding model $p(\mathbf{n} | x) = \int p(\mathbf{n} | s)p(s | x)ds$ contains information-limiting
345 correlations that bound the FI about x (*Moreno-Bote et al., 2014; Kanitscheider et al., 2015b*). To
346 understand whether and how the latent structure of CPMs captures information-limiting noise
347 correlations, we fit a minimal, von Mises, vanilla CPM with $d_K = 20$ mixture components to $d_T =$
348 10,000 responses from $p(\mathbf{n} | x)$. Figure 5A (purple) shows that the FI of the learned CPM saturates
349 near the precision of the sensory noise, indicating that the learned CPM accurately captures the
350 information-limiting correlations present in $p(\mathbf{n} | x)$.

351 To understand how the learned CPM represents the correlations in $p(\mathbf{n} | x)$ we study the re-
352 lation between the latent modulations and the population activity. Figure 5B shows the index-
353 probabilities of the learned CPM: given the true orientation $x = 90^\circ$, there are 3 modulations with
354 probabilities substantially greater than 0. To provide a high-level picture of how these modulations
355 affect population responses, in Figure 5C we plot the sum of the modulated rates of the population
356 as a function of orientation, and see that each modulation concentrates the tuning of the popula-
357 tion around a particular orientation, and that two of the modulations in particular shift the tuning
358 away from the true orientation.

359 Because there are essentially three modulations that are relevant to the responses of the CPM
360 to the true orientation $x = 90^\circ$, generating a response from the CPM approximately reduces to
361 generating a response from one of the three possible modulated populations. In Figures 5D-F
362 we depict a response to $x = 90^\circ$ from each of the three modulated populations, as well as the
363 optimal posterior based on the learned CPM (purple lines), and a suboptimal posterior based on
364 the source model (i.e. ignoring noise correlations; green lines). We observe that the trial-to-trial
365 variability of the learned CPM results in random shifts of the peak neural activity away from the
366 true orientation, thus fundamentally limiting information. Furthermore, when the response of the
367 population is concentrated at the true orientation (Figure 5E), the suboptimal posterior assigns a
368 high probability to the true orientation, whereas when the responses are biased away from the
369 true orientation (Figures 5D and 5F) the suboptimal posterior assigns nearly 0 probability to the
370 true orientation. This is in contrast to the optimal posterior, which always assigns a significant
371 probability to the true orientation.

372 In summary, CPMs accurately capture information-limiting correlations, and provide insight
373 into how such correlations can be generated by a simple latent structure.

374 Discussion

375 In this paper we introduced a latent variable exponential family formulation of multivariate Poisson
376 mixtures. We showed how this formulation allows us to effectively extend multivariate Poisson
377 mixtures both to capture sub-Poisson variability, and to incorporate stimulus dependence, which
378 we termed Conditional Poisson Mixtures (CPMs). Our analyses and simulations showed that CPMs
379 can be fit efficiently and recover ground truth models in synthetic data, capture a wide range of V1
380 response statistics in real data, and can be easily inverted to obtain accurate Bayesian decoding
381 that is competitive with nonlinear decoders, while using orders of magnitude less parameters. In
382 addition, we illustrated how the latent structure of CPMs provides an interpretable representation
383 of a fundamental feature of the neural code, e.g. information-limiting correlations.

384 Our framework is particularly relevant for probabilistic theories of neural coding based on the
385 theory of exponential families (*Beck et al., 2007*), which include theories that address the linearity
386 of Bayesian inference in neural circuits (*Ma et al., 2006*), the role of phenomena such as divisive
387 normalization in neural computation (*Beck et al., 2011a*), Bayesian inference about dynamic stim-
388 uli (*Makin et al., 2015; Sokolowski, 2017*), and the metabolic efficiency of neural coding (*Ganguli and*
389 *Simoncelli, 2014; Yerxa et al., 2020*). These theories have proven difficult to validate quantitatively
390 with neural data due to a lack of statistical models which are both compatible with their exponen-
391 tial family formulation, and can model correlated activity in recordings of large neural populations.
392 Our work suggests that CPMs can overcome these difficulties, and help connect the rich mathe-
393 matical theory of neural coding with the state-of-the-art in parallel recording technologies.

394 CPMs are not limited to modelling neural responses to stimuli, and can model how arbitrary
395 experimental variables modulate neural variability and covariability. Examples of experimental
396 variables that have measurable effects on neural covariability include the spatial and temporal
397 context around a stimulus (*Snyder et al., 2014; Snow et al., 2016, 2017; Festa et al., 2020*), as well
398 as task-variables and the attentional state of the animal (*Maunsell, 2015; Rabinowitz et al., 2015;*
399 *Kanashiro et al., 2017; Bondy et al., 2018; Ruff and Cohen, 2019*). Each of these variables could be
400 incorporated into a CPM by either replacing the stimulus-variable in our equations, or combining
401 it with the stimulus-variable to construct a CPM with multivariate dependence. This would allow
402 researchers to explore how the stimulus and the experimental variables mutually interact to shape
403 variability and covariability in large populations of neurons.

404 To understand how this variability and covariability effects neural coding, latent variable models
405 such as CPMs are often applied to extract interpretable features of the neural code from data (*White-*
406 *way and Butts, 2019*). The latent states of a CPM provide a soft classification of neural activity, and
407 we may apply CPMs to model how an experimental variable modulates the class membership of
408 neurons. In the aforementioned studies, models of neural activity yielded predictions of percep-
409 tual and behavioural performance. Because CPMs support Bayesian decoding, an appropriate
410 CPM can also make predictions about how a class of neurons is likely to modulate perception and
411 behaviour, and we may then test these predictions with experimental interventions on the neu-
412 rons themselves (*Panzeri et al., 2017*). In this manner, we believe CPMs could form a critical part
413 of a rigorous, Bayesian framework for “cracking the neural code” in large populations neurons.

414 In our applications we considered low-dimensional variables, and implemented the stimulus-
415 dependence of the CPM parameters with linear functions. Nevertheless, the stimulus-dependence
416 of a CPM can be implemented by arbitrary parametric functions of high-dimensional variables such
417 as deep neural networks, and CPMs can also incorporate history-dependence via recurrent neu-
418 ral networks. As such, CPMs have the potential to integrate encoding models of higher cortical
419 areas (*Yamins et al., 2014*) with models of the temporal features of the neural code (*Pillow et al.,*
420 *2008; Park et al., 2014; Runyan et al., 2017*), towards analyzing the neural code in dynamic, corre-
421 lated neural populations in higher cortex. Outside of neuroscience, high-dimensional count data
422 exists in many fields such as corpus linguistics and genomics (*Inouye et al., 2017*), and researchers
423 who aim to understand how this data depends on history or additional variables could benefit from
424 our techniques.

425 **Materials and methods**

426 **Notation**

427 We use capital, bold letters (e.g. Θ) to indicate matrices; small, bold letters (e.g. θ) to indicate vec-
428 tors; and regular letters (e.g. θ) to indicate scalars. We use subscript capital letters to indicate the
429 role of a given variable, so that, in Relation 1 for example, θ_K are the natural parameters that bias
430 the index-probabilities, θ_N are the baseline natural parameters of the neural firing rates, and Θ_{NK}
431 is the matrix of parameters through which the indices and rates interact.

432 We denote the i th element of a vector θ by θ_i , or e.g. of the vector θ_K by $\theta_{K,i}$. We denote the
433 i th row or j th column of Θ by θ_i or θ_j , respectively, and always state whether we are considering
434 a row or column of the given matrix. When referring to the j th element of a vector θ_i indexed by
435 i , we write θ_{ij} . Finally, when indexing data points from a sample, or parameters that are tied to
436 individual data points, we use parenthesized, superscript letters, e.g. $x^{(i)}$, or $\theta_N^{(i)}$.

437 **Poisson mixtures and their moments**

438 The following derivations were presented in a more general form in *Karlis and Meligkotsidou*
439 *(2007)*, but we present the simpler case here for completeness. A Poisson distribution has the form
440 $p(n; \lambda) = \frac{\lambda^n e^{-\lambda}}{n!}$, where n is the count and λ is the rate (in our case, spike count and firing rate, respec-
441 tively). We may use a Poisson model to define a distribution over d_N spike counts $\mathbf{n} = (n_1, \dots, n_{d_N})$

442 by supposing that the neurons generate spikes independently of one another, leading to the in-
 443 dependent Poisson model $p(\mathbf{n}; \lambda) = \prod_{i=1}^{d_N} p(n_i; \lambda_i)$ with firing rates $\lambda = (\lambda_1, \dots, \lambda_{d_N})$. Finally, if we
 444 consider the d_K rate vectors $\lambda_1, \dots, \lambda_{d_K}$, and d_K weights w_1, \dots, w_{d_K} , where $0 \leq w_k$ for all k , and
 445 $w_1 = 1 - \sum_{k=2}^{d_K} w_k$, we then define a mixture of Poisson distributions as a latent variable model
 446 $p(n) = \sum_k p(n | k)p(k) = \sum_k p(n, k)$, where $p(n | k) = p(n; \lambda_k)$, and $p(k) = w_k$.

447 The mean μ_i of the i th neuron of a mixture of independent Poisson distributions is

$$\mu_i = \sum_{n_i=0}^{\infty} \sum_{k=1}^{d_K} p(n_i | k)p(k)n_i = \sum_{k=1}^{d_K} p(k) \sum_{n_i=0}^{\infty} p(n_i | k)n_i = \sum_{k=1}^{d_K} w_k \lambda_{ik}. \quad (4)$$

448 The variance σ_i^2 of neuron i is

$$\sigma_i^2 = \sum_{n_i=0}^{\infty} p(n_i)n_i^2 - \mu_i^2 = \sum_{k=1}^{d_K} p(k) \sum_{n_i=0}^{\infty} p(n_i | k)n_i^2 - \mu_i^2 = \sum_{k=1}^{d_K} p(k)(\sigma_{ik}^2 + \lambda_{ik}^2) - \mu_i^2 = \mu_i + \sum_{k=1}^{d_K} w_k(\lambda_{ik} - \mu_i)^2, \quad (5)$$

449 where $\sigma_{ik}^2 = \lambda_{ik}$ is the variance of the i th neuron under the k th component distribution, i.e. the
 450 variance of $p(n_i | k)$, and where $\sum_{n_i=0}^{\infty} p(n_i | k)n_i^2 = \sigma_{ik}^2 + \lambda_{ik}^2$, and $\sum_{k=1}^{d_K} w_k \lambda_{ik}^2 - \mu_i^2 = \sum_{k=1}^{d_K} w_k(\lambda_{ik} - \mu_i)^2$
 451 both follow from the fact that a distribution's variance equals the difference between its second
 452 moment and squared first moment.

The covariance σ_{ij}^2 between spike-counts n_i and n_j for $i \neq j$ is then

$$\begin{aligned} \sigma_{ij}^2 &= \sum_{n_i=0}^{\infty} \sum_{n_j=0}^{\infty} p(n_i, n_j)(n_i - \mu_i)(n_j - \mu_j) = \sum_{k=1}^{d_K} p(k) \sum_{n_i=0}^{\infty} \sum_{n_j=0}^{\infty} p(n_i, n_j | k)(n_i - \mu_i)(n_j - \mu_j) \\ &= \sum_{k=1}^{d_K} p(k) \sum_{n_i=0}^{\infty} p(n_i | k)(n_i - \mu_i) \sum_{n_j=0}^{\infty} p(n_j | k)(n_j - \mu_j) = \sum_{k=1}^{d_K} w_k(\lambda_{ik} - \mu_i)(\lambda_{jk} - \mu_j). \end{aligned} \quad (6)$$

453 Observe that if $w_k = \frac{1}{d_K - 1}$, then σ_{ij}^2 is simply the sample covariance between i and j , where the
 454 sample is composed of the rate components of the i th and j th neurons. Equation 6 thus implies
 455 that Poisson mixtures can model arbitrary covariances. Nevertheless, Equation 5 shows that the
 456 variance of individual neurons is restricted to being larger than their means.

457 Exponential family mixture models

458 In this section we show that the latent variable form for Poisson mixtures we introduced above
 459 is a member of the class of models known as exponential families. An exponential family distri-
 460 bution $p(x)$ over some data x has the form $p(x) = e^{\theta \cdot \mathbf{s}(x) - \psi(\theta)} b(x)$, where θ are the so-called nat-
 461 ural parameters, $\mathbf{s}(x)$ is a vector-valued function of the data called the sufficient statistic, $b(x)$ is
 462 a scalar-valued function called the base measure, and $\psi(\theta) = \log \int e^{\theta \cdot \mathbf{s}(x)} b(x) dx$ is the log-partition
 463 function (*Wainwright and Jordan, 2008*). In the context of Poisson mixture models, we note that an
 464 independent Poisson model $p(\mathbf{n}; \lambda)$ is an exponential family, with natural parameters θ_N given by
 465 $\theta_{N,i} = \log \lambda_i$, base measure $b(\mathbf{n}) = \prod_i n_i!$ and sufficient statistic $\mathbf{s}_N(\mathbf{n}) = \mathbf{n}$, and log-partition function
 466 $\psi_N(\theta_N) = \log \sum_i e^{\theta_{N,i}}$. Moreover, the distribution of component indices $p(k)$ (also known as a cat-
 467 egorical distribution) also has an exponential family form, with natural parameters $\theta_{K,k} = \log \frac{w_{k+1}}{w_1}$
 468 for $1 \leq k < d_K$, sufficient statistic $\delta(k) = (\delta_2(k), \dots, \delta_{d_K}(k))$, base measure $b(k) = 1$, and log-partition
 469 function $\psi_K(\theta_K) = \log(1 + \sum_{k=1}^{d_K-1} e^{\theta_{K,k}})$. Note that in both cases, the exponential parameters are well-
 470 defined only if the rates and weights are strictly greater than 0 — in practice however this is not a
 471 significant limitation.

472 We claim that the joint distribution of a multivariate Poisson mixture model $p(\mathbf{n}, k)$ can be repa-
 473 rameterized in the exponential family form

$$p(\mathbf{n}, k) = \frac{e^{\theta_N \cdot \mathbf{n} + \theta_K \cdot \delta(k) + \mathbf{n} \cdot \theta_{NK} \cdot \delta(k) - \psi_{NK}(\theta_N, \theta_K, \theta_{NK})}}{\prod_i n_i!}, \quad (7)$$

474 where $\psi_{NK}(\theta_N, \theta_K, \Theta_{NK}) = \log \sum_{\mathbf{k}} e^{\theta_K \cdot \delta(k) + \psi_N(\theta_N + \Theta_{NK} \cdot \delta(k))}$ is the log-partition function of $p(\mathbf{n} | k)$. To
 475 show this we show how to express the natural parameters θ_N , θ_K , and Θ_{NK} as (invertible) functions
 476 of the component rate vectors $\lambda_1, \dots, \lambda_{d_K}$, and the weights w_1, \dots, w_{d_K} . In particular, we set

$$\theta_N = \log \lambda_1, \quad (8)$$

477 where log is applied element-wise. Then, for $1 \leq k < d_K$, we set the k th row $\theta_{NK,k}$ of Θ_{NK} to

$$\theta_{NK,k} = \log \lambda_{k+1} - \log \lambda_1, \quad (9)$$

478 and the k th element of θ_K to

$$\theta_{K,k} = \log \frac{w_{k+1}}{w_1} + \psi(\theta_N) - \psi_N(\theta_N + \Theta_{NK} \cdot \delta(k)). \quad (10)$$

479 This reparameterization may then be checked by substituting Equations 8, 9, and 10 into Equation 7
 480 to recover the joint distribution of the mixture model $p(\mathbf{n}, k) = p(\mathbf{n} | k)p(k) = w_k p(\mathbf{n}; \lambda_K)$; for a more
 481 explicit derivation see [Sokoloski \(2019\)](#).

482 The equation for $p(\mathbf{n}, k)$ ensures that the index-probabilities are given by

$$p(k) = e^{\theta_K \cdot \delta(k) - \psi_{NK}(\theta_N, \theta_K, \Theta_{NK})} \sum_{\mathbf{n}} \frac{e^{\mathbf{n} \cdot (\theta_N + \Theta_{NK} \cdot \delta(k))}}{\prod_i n_i!} = e^{\theta_K \cdot \delta(k) - \psi_{NK}(\theta_N, \theta_K, \Theta_{NK}) + \psi_N(\theta_N + \Theta_{NK} \cdot \delta(k))}. \quad (11)$$

483 Consequently, the component distributions in exponential family form are given by

$$p(\mathbf{n} | k) = \frac{p(\mathbf{n}, k)}{p(k)} = e^{\mathbf{n} \cdot (\theta_N + \Theta_{NK} \cdot \delta(k)) - \psi_N(\theta_N + \Theta_{NK} \cdot \delta(k))}. \quad (12)$$

484 Observe that $p(\mathbf{n} | k)$ is a multivariate Poisson distribution with parameters $\theta_N + \Theta_{NK} \cdot \delta(k)$, so
 485 that for $k > 1$, the parameters are the sum of θ_N and row $k - 1$ of Θ_{NK} . Because the exponential
 486 family parameters are the logarithms of the firing rates of \mathbf{n} , each row of Θ_{NK} modulates the firing
 487 rates of \mathbf{n} multiplicatively. When $\theta_N(x)$ depends on a stimulus and we consider the component
 488 distributions $p(\mathbf{n} | x, k)$, each row of Θ_{NK} then scales the tuning curves of the baseline population
 489 (i.e. $p(\mathbf{n} | x, k)$ for $k = 1$); in the neuroscience literature, such scaling factors are typically referred
 490 to as gain modulations.

491 The exponential family form has many advantages. However, it has a less intuitive relationship
 492 with the statistics of the model such as the mean and covariance. The most straightforward method
 493 to compute these statistics given a model in exponential family form is to first reparameterize it in
 494 terms of the weights and component rates, and then evaluate Equations 4, 5, and 6.

495 CoM-Poisson distributions and their mixtures

496 Conway-Maxwell (CoM) Poisson distributions decouple the location and shape of count distribu-
 497 tions ([Shmueli et al., 2005](#); [Stevenson, 2016](#); [Chanialidis et al., 2018](#)). A CoM Poisson model has
 498 the form $p(n; \lambda, \nu) \propto \left(\frac{\lambda^n}{n!}\right)^\nu$. The floor function $\lfloor \lambda \rfloor$ of the location parameter λ is the mode of the given
 499 distribution. With regards to the shape parameter ν , $p(n; \lambda, \nu)$ is a Poisson distribution with rate λ
 500 when $\nu = 1$, and is under- or over-dispersed when $\nu > 1$ or $\nu < 1$, respectively. A CoM-Poisson model
 501 $p(n; \lambda, \nu)$ is also an exponential family, with natural parameters $\theta_C = (\nu \log \lambda, -\nu)$, sufficient statistic
 502 $s_C(n) = (n, \log n!)$, and base measure $b(n) = 1$. The log-partition function does not have a closed-form
 503 expression, but it can be effectively approximated by truncating the series $\sum_{n=0}^{\infty} e^{s_C(n) \cdot \theta_C}$ ([Shmueli](#)
 504 [et al., 2005](#)). More generally, when we consider a product of independent CoM-Poisson distri-
 505 butions, we denote its log-partition function by $\log \psi_C(\theta_N, \theta_N^*) = \sum_{i=1}^{d_N} \sum_{n_i=0}^{\infty} e^{n_i \theta_{N,i} + \log(n_i)! \theta_{N,i}^*}$, where
 506 $\theta_{C,i} = (\theta_{N,i}, \theta_{N,i}^*)$ are the parameters of the i th CoM-Poisson distribution. In this case we can also ap-
 507 proximate the log-partition function ψ_C by truncating the d_N constituent series $\sum_{n_i=0}^{\infty} e^{n_i \theta_{N,i} + \log(n_i)! \theta_{N,i}^*}$
 508 in parallel.

509 We define a multivariate CoM-based mixture as

$$p(\mathbf{n}, k) = e^{\theta_N \cdot \mathbf{n} + \theta_N^* \cdot \mathbf{lf}(\mathbf{n}) + \theta_K \cdot \delta(k) + \mathbf{n} \cdot \Theta_{NK} \cdot \delta(k) - \psi_{CK}(\theta_N, \theta_N^*, \theta_K, \Theta_{NK})}, \quad (13)$$

510 where $\mathbf{f}(\mathbf{n}) = (\log(n_1!), \dots, \log(n_{d_N}!))$ is the vector of log-factorials of the individual spike-counts,
 511 and $\psi_{CK}(\theta_N, \theta_N^*, \theta_K, \Theta_{NK}) = \log \sum_k e^{\theta_k \cdot \delta(k) + \psi_C(\theta_N + \Theta_{NK} \cdot \delta(k), \theta_N^*)}$ is the log-partition function. This form
 512 ensures that the index-probabilities satisfy

$$p(k) = e^{\theta_k \cdot \delta(k) - \psi_{CK}(\theta_N, \theta_N^*, \theta_K, \Theta_{NK}) + \psi_C(\theta_N + \Theta_{NK} \cdot \delta(k), \theta_N^*)}, \quad (14)$$

513 and consequently that each component distribution $p(\mathbf{n} | k)$ is a product of independent CoM
 514 Poisson distributions given by

$$p(\mathbf{n} | k) = e^{\mathbf{n} \cdot (\theta_N + \Theta_{NK} \cdot \delta(k)) + \theta_N^* \cdot \mathbf{f}(\mathbf{n}) - \psi_C(\theta_N + \Theta_{NK} \cdot \delta(k), \theta_N^*)}. \quad (15)$$

515 Observe that, whereas the parameters $\theta_N + \Theta_{NK} \cdot \delta(k)$ of $p(\mathbf{n} | k)$ depend on the index k , the
 516 parameters θ_N^* of $p(\mathbf{n} | k)$ are independent of the index and act exclusively as biases. Note as well
 517 that when considering a CoM-based, minimal CPM, the modulated populations ($p(\mathbf{n} | k, x)$ for $k > 1$)
 518 continue to scale the firing rates of the baseline population ($p(\mathbf{n} | k, x)$) monotonically, but not in a
 519 linear, multiplicative manner.

520 The moments of a CoM-Poisson distribution are not available in closed-form, yet they can also
 521 be effectively approximated through truncation. Given approximate means μ_{ik} and variances σ_{ik}^2
 522 of $p(n_i | k)$, we may easily evaluate the means, variances, and covariances of $p(n_i)$. In particular, the
 523 mean of n_i is $\mu_i = \sum_{k=1}^{d_K} p(k) \mu_{ik}$, and its variance is

$$\sigma_i^2 = \bar{\sigma}_i^2 + \sum_{k=1}^{d_K} p(k) (\mu_{ik} - \mu_i)^2, \quad (16)$$

524 where $\bar{\sigma}_i^2 = \sum_{k=1}^{d_K} p(k) \sigma_{ik}^2$. Finally, similarly to Equation 6, the covariance σ_{ij} between n_i and n_j is
 525 $\sigma_{ij} = \sum_{k=1}^{d_K} p(k) (\mu_{ik} - \mu_i) (\mu_{jk} - \mu_j)$.

526 By comparing Equations 5 and 16, we see that the CoM-based mixture may address the lim-
 527 itations on the variances σ_i^2 of the vanilla mixture by setting the average variance $\bar{\sigma}_i^2$ of the com-
 528 ponents in Equation 16 to be small, while holding the value of the means μ_i fixed, and ensuring
 529 that the means of the components μ_{ik} cover a wide range of values to achieve the desired values
 530 of σ_i^2 and σ_{ij} . Solving the parameters of a CoM-based mixture for a desired covariance matrix is
 531 unfortunately not possible since we lack closed-form expressions for the means and variances.
 532 Nevertheless, we may justify the effectiveness of the CoM-based strategy by considering the ap-
 533 proximations of the components means and variances $\mu_{ik} \approx \lambda_{ik} + \frac{1}{2v_{ik}} - \frac{1}{2}$ and $\sigma_{ik}^2 \approx \frac{\lambda_{ik}}{v_{ik}}$, which hold
 534 when neither λ_{ik} or v_{ik} are too small (*Chaniyalidis et al., 2018*). Based on these approximations,
 535 observe that when v_{ik} is large, σ_{ik}^2 is small, whereas μ_{ik} is more or less unaffected. Therefore, in
 536 the regime where these approximations hold, a small value for $\bar{\sigma}_i^2$ can be achieved by reducing the
 537 parameters v_{ik} , without significantly restricting the values of μ_{ik} or μ_i .

538 Fisher information of a CPM

The Fisher information (FI) of an encoding model $p(\mathbf{n} | x)$ with respect to x is $I(x) = \sum_{\mathbf{n}} p(\mathbf{n} | x) (\partial_x \log p(\mathbf{n} | x))^2$ (*Cover and Thomas, 2006*). With regards to the FI of a CPM,

$$\begin{aligned} \partial_x \log p(\mathbf{n} | x) &= \frac{\sum_k \partial_x p(\mathbf{n}, k | x)}{p(\mathbf{n} | x)} = \frac{\sum_k \partial_x e^{\theta_N(x) \cdot \mathbf{n} + \theta_N^* \cdot \mathbf{f}(\mathbf{n}) + \theta_K \cdot \delta(k) + \mathbf{n} \cdot \Theta_{NK} \cdot \delta(k) - \psi_{CK}(\theta_N(x), \theta_N^*, \theta_K, \Theta_{NK})}}{p(\mathbf{n} | x)} \\ &= \partial_x (\theta_N(x) \cdot \mathbf{n} - \psi_{CK}(\theta_N(x), \theta_N^*, \theta_K, \Theta_{NK})) \frac{\sum_k p(\mathbf{n}, k | x)}{p(\mathbf{n} | x)} = \partial_x \theta_N(x) \cdot (\mathbf{n} - \mu_N(x)), \end{aligned}$$

where $\partial_x \psi_{CK}(\theta_N(x), \theta_N^*, \theta_K, \Theta_{NK}) = \mu_N(x) \cdot \partial_x \theta_N(x)$ follows from the chain rule and properties of the log-partition function (*Wainwright and Jordan, 2008*). Therefore

$$I(x) = \sum_{\mathbf{n}} p(\mathbf{n} | x) (\partial_x \theta_N(x) \cdot (\mathbf{n} - \mu_N(x)))^2 = \partial_x \theta_N(x) \cdot \Sigma_N(x) \cdot \partial_x \theta_N(x),$$

539 where $\Sigma_N(x)$ is the covariance matrix of $p(\mathbf{n} | x)$. Moreover, because $\partial_x \theta_N(x) = \Sigma_N^{-1}(x) \cdot \partial_x \mu(x)$ (*Wain-*
 540 *wright and Jordan, 2008*), the FI of a CPM may also be expressed as $I(x) = \partial_x \mu_N(x) \cdot \Sigma_N^{-1}(x) \cdot \partial_x \mu_N(x)$,
 541 which is the linear Fisher information (*Beck et al., 2011b*).

542 Note that when calculating the FI or other quantities based on the covariance matrix, vanilla
 543 CPMs have the advantage that their covariance matrices tend to have large diagonal elements and
 544 are thus inherently well-conditioned. Because decoding performance is not significantly different
 545 between vanilla and CoM-based CPMs (see Table 2), vanilla CPMs may be preferable when well-
 546 conditioned covariance matrices are critical. Nevertheless, the covariance matrices of CoM-based
 547 mixtures can be made well-conditioned by applying standard techniques.

548 Expectation-Maximization for CPMs

549 Expectation-maximization (EM) is an algorithm that maximizes the likelihood of a latent variable
 550 model given data by iterating two steps: generating model-based expectations of the latent vari-
 551 ables, and maximizing the complete log-likelihood of the model given the data and latent expecta-
 552 tions. Although the maximization step optimizes the *complete* log-likelihood, each iteration of EM
 553 is guaranteed to increase the *data* log-likelihood as well (*Neal and Hinton, 1998*).

554 EM is arguably the most widely-applied algorithm for fitting finite mixture models (*McLachlan*
 555 *et al., 2019*). As a form of latent variable exponential family, the expectation step for a finite mixture
 556 model reduces to computing average sufficient statistics, and the maximization step is a convex
 557 optimization problem (*Wainwright and Jordan, 2008*). In general, the average sufficient statistics,
 558 or mean parameters, correspond to (are dual to) the natural parameters of an exponential family,
 559 and where we denote natural parameters with θ , we denote their corresponding mean parameters
 560 with η .

Suppose we are given a dataset $(\mathbf{n}^{(1)}, \dots, \mathbf{n}^{(d_T)})$ of neural spike-counts, and a CoM-based mixture
 model with natural parameters $\theta_N, \theta_N^*, \theta_K$, and Θ_{NK} (see Equation 13). The expectation step for
 this model reduces to computing the data-dependent mean parameters $\eta_K^{(i)}$ given by

$$\theta_K^{(i)} = \theta_K + \mathbf{n}^{(i)} \cdot \Theta_{NK}, \quad \eta_{K,k}^{(i)} = \frac{e^{\theta_{K,k}^{(i)}}}{1 + \sum_l e^{\theta_{K,l}^{(i)}}},$$

561 for all $0 < i \leq d_T$. The mean parameters $\eta_K^{(i)}$ are the averages of the sufficient statistic $\delta_k(k)$ under
 562 the distribution $p(k | \mathbf{n}^{(i)})$, and are what we use to complete the log-likelihood since we do not
 563 observe k .

Given $\eta_K^{(i)}$, the maximization step of a CoM-based mixture thus reduces to maximizing the com-
 plete log-likelihood $\sum_{i=1}^{d_T} \mathcal{L}(\theta_K, \theta_N, \theta_N^*, \Theta_{NK}, \eta_K^{(i)}, \mathbf{n}^{(i)})$, where we substitute $\eta_K^{(i)}$ into the place of $\delta(k)$
 in Equation 13, such that

$$\mathcal{L}(\theta_K, \theta_N, \theta_N^*, \Theta_{NK}, \eta_K^{(i)}, \mathbf{n}^{(i)}) = \theta_N \cdot \mathbf{n}^{(i)} + \theta_N^* \cdot \mathbf{I}(\mathbf{n}^{(i)}) + \theta_K \cdot \eta_K^{(i)} + \mathbf{n}^{(i)} \cdot \Theta_{NK} \cdot \eta_K^{(i)} - \psi_{CK}(\theta_N, \theta_N^*, \theta_K, \Theta_{NK}).$$

564 This objective may be maximized in closed-form for a vanilla Poisson mixture (*Karlis and Meligkoti-*
 565 *sidou, 2007*), but this is not the case when the model has CoM-Poisson shape parameters or de-
 566 pends on the stimulus. Nevertheless, solving the resulting maximization step is still a convex opti-
 567 mization problem (*Wainwright and Jordan, 2008*), and may be approximately solved with gradient
 568 ascent. Doing so requires that we first compute the mean parameters η_N, η_N^*, η_K , and \mathbf{H}_{NK} that
 569 are dual to $\theta_N, \theta_N^*, \theta_K$, and Θ_{NK} , respectively.

We compute the mean parameters by evaluating

$$\theta_{K,k}^\dagger = \theta_{K,k} + \psi_C(\theta_N + \Theta_{NK} \cdot \delta(k), \theta_N^*) - \psi(\theta_N), \quad \eta_{K,k} = \frac{e^{\theta_{K,k}^\dagger}}{1 + \sum_{k=1}^{d_K-1} e^{\theta_{K,k}^\dagger}}, \quad \mu_{jk} = \sum_{n_j=0}^{\infty} n_j p(n_j | k),$$

$$\eta_{N,j}^* = \sum_{k=1}^{d_K} p(k) \sum_{n_j=0}^{\infty} \log n_j! p(n_j | k), \quad \eta_{N,j} = \sum_{k=1}^{d_K} p(k) \mu_{jk}, \quad \eta_{NK,jk} = \eta_{K,k} \mu_{j(k+1)},$$

where $\eta_{K,k}$ is the k th element of η_K , $\eta_{N,j}$ is the j th element of η_N , $\eta_{N,j}^*$ is the j th element of η_N^* ,
 and $\eta_{NK,jk}$ is the j th element of the k th column of \mathbf{H}_{NK} . Note as well that we truncate the series

$\sum_{n_j} n_j p(n_j | k)$ and $\sum_{n_j} \log n_j! p(n_j | k)$ to approximate μ_{jk} and $\eta_{N,j}^*$. Given these mean parameters, we may then express the gradients of $\mathcal{L}^{(i)} = \mathcal{L}(\theta_K, \theta_N, \theta_N^*, \Theta_{NK}, \eta_{K,i}, \mathbf{n}^{(i)})$ as

$$\begin{aligned} \partial_{\theta_N} \mathcal{L}^{(i)} &= \mathbf{n}^{(i)} - \eta_N, & \partial_{\theta_N^*} \mathcal{L}^{(i)} &= \mathbf{I} \mathbf{n}^{(i)} - \eta_N^*, \\ \partial_{\theta_K} \mathcal{L}^{(i)} &= \eta_K^{(i)} - \eta_K, & \partial_{\Theta_{NK}} \mathcal{L}^{(i)} &= \mathbf{n}^{(i)} \otimes \eta_K^{(i)} - \mathbf{H}_{NK}, \end{aligned}$$

570 where \otimes is the outer product operator, and where the second term in each equation follows from
571 the fact that the derivative of ψ_{CK} with respect to θ_N , θ_N^* , θ_K , or Θ_{NK} yields the dual parameters
572 η_N , η_N^* , η_K , and \mathbf{H}_{NK} , respectively. By ascending the gradients of $\sum_{i=1}^{d_T} \mathcal{L}^{(i)}$ until convergence, we
573 approximate a single iteration of the EM algorithm for a CoM-based mixture.

Finally, if our dataset $((\mathbf{n}^{(1)}, x^{(1)}), \dots, (\mathbf{n}^{(d_T)}, x^{(d_T)}))$ includes stimuli x , and the parameters θ_N depend on the stimulus, then the gradients of the parameters of θ_N must also be computed. For a von Mises CPM where $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \mathbf{vm}(x)$, the gradients are given by

$$\partial_{\theta_N^0} \mathcal{L}^{(i)} = \partial_{\theta_N^{(i)}} \mathcal{L}^{(i)}, \quad \partial_{\Theta_{NX}} \mathcal{L}^{(i)} = \partial_{\theta_N^{(i)}} \mathcal{L}^{(i)} \otimes \mathbf{vm}(x^{(i)}),$$

574 where $\theta_N^{(i)} = \theta_N(x^{(i)})$ is the output of θ_N at $x^{(i)}$. Although in this paper we restrict our applications
575 to Von Mises or discrete tuning curves for 1-dimensional stimuli, this formalism can be readily
576 extended to the case where the baseline tuning curve parameters $\theta_N(x)$ are a generic nonlinear
577 function of the stimulus, represented by a deep neural network. Then, the gradients of the pa-
578 rameters of θ_N can be computed through backpropagation, and $\partial_{\theta_N^{(i)}} \mathcal{L}^{(i)}$ is the error that must be
579 backpropagated through the network to compute the gradients.

580 CPM initialization and training procedures

581 To fit a CPM to a dataset $((\mathbf{n}^{(1)}, x^{(1)}), \dots, (\mathbf{n}^{(d_T)}, x^{(d_T)}))$, we first initialize the CPM and then optimize
582 its parameters with our previously described EM algorithm. Naturally, initialization depends on
583 exactly which form of CPM we consider, but in general we first initialize the baseline parameters θ_N ,
584 then add the categorical parameters θ_K and mixture component parameters Θ_{NK} . When training
585 CoM-based CPMs we always first train a vanilla CPM, and so the initialization procedure remains
586 the same for vanilla and CoM-based models.

To initialize a minimal, von Mises CPM with d_N neurons, we first fit d_N independent, von Mises-
tuned neurons by maximizing the log-likelihood $\sum_{i=1}^{d_T} \log p(\mathbf{n}^{(i)} | x^{(i)})$ of $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \mathbf{vm}(x)$.
This is a convex optimization problem and so can be easily solved by gradient ascent, in particular
by following the gradients

$$\begin{aligned} \partial_{\theta_N^0} \sum_{i=1}^{d_T} \log p(\mathbf{n}^{(i)} | x^{(i)}) &= \sum_{i=1}^{d_T} \mathbf{n}^{(i)} - \log(\theta_N(x^{(i)})), \\ \partial_{\Theta_{NX}} \sum_{i=1}^{d_T} \log p(\mathbf{n}^{(i)} | x^{(i)}) &= \sum_{i=1}^{d_T} \log(\mathbf{n}^{(i)} - \log \theta_N(x^{(i)})) \otimes \mathbf{vm}(x^{(i)}), \end{aligned}$$

587 to convergence. For both discrete and maximal CPMs, where there are d_X distinct stimuli, we
588 initialize $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \delta(x)$ by computing the average rate vector at each stimulus-condition
589 and creating a lookup table for these rate vectors. Formally, where x_l is the l th stimulus value for
590 $0 < l \leq d_X$, we may express the l th rate vector as $\lambda_l = \frac{1}{\sum_{i=1}^{d_T} \delta(x_l, x^{(i)})} \sum_{i=1}^{d_T} \delta(x_l, x^{(i)}) \mathbf{n}^{(i)}$, where $\delta(x_l, x^{(i)})$
591 is 1 when $x_l = x^{(i)}$, and 0 otherwise. We then construct a lookup table for these rate vectors in
592 exponential family form by setting $\theta_N^0 = \log \lambda_1$, and by setting the l th row $\theta_{NX,l}$ of Θ_{NX} to $\theta_{NX,l} =$
593 $\log \lambda_{l+1} - \log \lambda_1$.

594 In general we initialize the parameters θ_K by sampling the weights w_1, \dots, w_{d_K} of a categori-
595 cal distribution from a Dirichlet distribution with a constant concentration of 2, and converting
596 the weights into the natural parameters of a categorical distribution θ_K . For discrete and maxi-
597 mal CPMs we initialize the modulations Θ_{NK} by generating each element of Θ_{NK} from a uniform
598 distribution over the range $[-0.0001, 0.0001]$. For von Mises CPMs we initialize each row $\theta_{NK,k}$ of

599 Θ_{NK} as shifted sinusoidal functions of the preferred stimuli of the independent von Mises neu-
 600 rons. That is, given θ_N^0 and Θ_{NX} , we compute the preferred stimulus of the i th neuron given by
 601 $\rho_i = \text{atan2}(\theta_N^0 + \Theta_{NX,i})$, where $\Theta_{NX,i}$ is the i th row of Θ_{NX} . We then set the i th element $\theta_{NK,k,i}$ of $\Theta_{NK,k}$
 602 to $\theta_{NK,k,i} = 0.2 \sin(\rho_i + \frac{k}{360}^\circ)$. Initializing von Mises CPMs in this way ensures that each modulation
 603 has a unique peak as a function of preferred stimuli, which helps differentiate the modulations
 604 from each other, and in our experience improves training speed.

With regards to training, the expectation step in our EM algorithm may be computed directly,
 and so the only challenge is solving the maximization step. Although the optimal solution strategy
 depends on the details of the model and data in question, in the context of this paper we settled on
 a strategy that is sufficient for all simulations we perform. For each model we perform a total of $d_I =$
 500 EM iterations, and for each maximization step we take $d_S = 100$ gradient ascent steps with the
 Adam gradient ascent algorithm (Kingma and Ba, 2014) with the default momentum parameters
 (see Kingma and Ba (2014)). We restart the Adam algorithm at each iteration of EM and gradually
 reduce the learning rate. Where $\epsilon^+ = 0.002$ and $\epsilon^- = 0.0005$ are the initial and final learning rates,
 we set the learning rate ϵ_t at EM iteration t to

$$\epsilon_t = \exp\left(\frac{(d_I - 1 - t) \log(\epsilon^+) + t \log(\epsilon^-)}{d_I - 1}\right),$$

605 where we assume t starts at 0 and ends at $d_I - 1$.

606 Because we must evaluate large numbers of truncated series when working with CoM-based
 607 CPMs, training times are typically one to two orders of magnitude greater. To minimize training
 608 time of CoM-based CPMs over the d_I EM iterations, we therefore first train a vanilla CPM for $0.8d_I$
 609 iterations. We then equate the parameters θ_N , θ_K , and Θ_{NK} of the vanilla CPM (see Equation 7)
 610 with a CoM-based CPM (see Equation 13) and set $\theta_N^* = -\mathbf{1}$, which ensures that resulting CoM-based
 611 model has the same density function $p(\mathbf{n}, k | x)$ as the original vanilla model. We then train the CoM-
 612 based CPM for $0.2d_I$ iterations. We found this strategy results in practically no performance loss,
 613 while greatly reducing training time.

614 CPM parameter selection for simulations

615 In the section Extended Poisson mixture models capture stimulus-dependent response statistics
 616 and the section CPMs facilitate accurate and efficient Bayesian decoding of neural responses we
 617 considered CoM-based, minimal CPMs with randomized parameters $\theta_N(x)$, θ_N^* , θ_K , and Θ_{NK} , which
 618 for simplicity we refer to as models 1 and 2, respectively. We construct randomized CPMs piece by
 619 piece, in a similar fashion to our initialization procedure.

620 Firstly, where d_N is the number of neurons, we tile their preferred stimuli ρ_i over the circle such
 621 that $\rho_i = \frac{i}{d_N} 360^\circ$. We then generate the concentration κ_i and gain γ_i of the i th neuron by sampling
 622 from normal distributions in log-space, such that $\log \kappa_i \sim N(-0.1, 0.2)$, and $\log \gamma_i \sim N(0.2, 0.1)$. Finally,
 623 for von Mises baseline tuning curves $\theta_N(x) = \theta_N^0 + \Theta_{NX} \cdot \mathbf{vm}(x)$, we set each row $\theta_{NX,i}$ of Θ_{NX} to
 624 $\theta_{NX,i} = (\kappa_i \cos \rho_i, \kappa_i \sin \rho_i)$, and each element $\theta_{N,i}^0$ of θ_N^0 to $\theta_{N,i}^0 = \log \gamma_i - \psi_X(\theta_{NX,i})$, where ψ_X is the
 625 logarithm of the modified Bessel function of order 0, which is the log-partition function of the von
 626 Mises distribution.

627 We then set $\theta_K = \mathbf{0}$, and generated each element $\theta_{NK,i,k}$ of the modulation matrix Θ_{NK} in the
 628 same matter as the gains, such that $\theta_{NK,i,k} \sim N(0.2, 0.1)$. Finally, to generate random CoM-based
 629 parameters we generate each element $\theta_{N,i}^*$ of θ_N^* from a uniform distribution, such that $\theta_{N,i}^* \sim$
 630 $U(-1.5, -0.8)$.

631 Model 2 entails two more steps. Firstly, when sampling from larger populations of neurons,
 632 single modulations often dominate the model activity around certain stimulus values. To suppress
 633 this we consider the natural parameters $\theta_{K,k}^0(x)$ of $p(k | x)$ (see Equation 14), and compute the max-
 634 imum value of these natural parameters over the range of stimuli $\theta_{K,k}^+ = \max_x \{\theta_{K,k}^0(x)\}$. We then
 635 set each element $\theta_{K,k}$ of the parameters θ_K of the CPM to $\theta_{K,k} = \bar{\theta}_K^+ - \theta_{K,k}^+$, where $\bar{\theta}_K^+ = \sum_{i=1}^{d_K} \frac{\theta_{K,k}}{d_K}$,
 636 which helps ensure that multiple modulations are active at any given x . Finally, since model 2

637 is a discrete CPM, we replace the von Mises baseline tuning curves with discrete baseline tuning
638 curves, by evaluating $\theta_N^0 + \Theta_{NX} \cdot \mathbf{vm}(x)$ at each of the d_x valid stimulus-conditions, and assemble
639 the resulting collection of natural parameters into a lookup table in the manner we described in
640 our initialization procedures.

641 Decoding models

642 When constructing a Bayesian decoder for discrete stimuli, we first estimate the prior $p(x)$ by com-
643 puting the relative frequency of stimulus presentations in the training data. For the given encod-
644 ing model, we then evaluate $p(\mathbf{n} | x)$ at each stimulus condition, and then compute the posterior
645 $p(x | \mathbf{n}) \propto p(\mathbf{n} | x)p(x)$ by brute-force normalization of $p(\mathbf{n} | x)p(x)$. When training the encoding
646 model used for our Bayesian encoders, we only trained them to maximize encoding performance
647 as previously described, and not to maximize decoding performance.

648 We considered two decoding models, namely the linear network and the artificial neural net-
649 work (ANN) with sigmoid activation functions. In both cases the input of the network was a neural
650 response vector, and the output the natural parameters θ_x of a categorical distribution. The form
651 of the linear network was $\theta_x(\mathbf{n}) = \theta_x + \Theta_{xN} \cdot \mathbf{n}$, and is otherwise fully determined by the structure
652 of the data. For the ANN on the other hand, we had to choose both the number of hidden layers,
653 and the number of neurons per hidden layer. We cross-validated the performance of both 1 and 2
654 hidden layer models, over a range of sizes from 100 to 2000 neurons. We found the performance
655 of the networks with 2 hidden layers generally exceeded that of those with 1 hidden layer, and that
656 700 and 600 hidden neurons was optimal for the awake and anaesthetized networks, respectively.

657 Given a dataset $((\mathbf{n}^{(1)}, x^{(1)}), \dots, (\mathbf{n}^{(d_T)}, x^{(d_T)}))$, we optimized the linear network and the ANN by max-
658 imizing $\sum_{i=1}^{d_T} \log p(x^{(i)} | \mathbf{n}^{(i)})$ via stochastic gradient ascent. We again used the Adam optimizer with
659 default momentum parameters, and used a fixed learning rate of 0.0003, and randomly divided
660 the dataset into minibatches of 500 data points. We also used early stopping, where for each fold
661 of our 10-fold cross-validation simulation, we partitioned the dataset into 80% training data, 10%
662 test data, and 10% validation data, and stopped the simulation when performance on the test data
663 declined from epoch to epoch.

664 Experimental design

665 Throughout this paper we demonstrate our methods on two sets of parallel response recordings
666 in macaque primary visual cortex (V1). The stimuli were drifting full contrast gratings at 9 distinct
667 orientations spread evenly over the half-circle from 0° to 180° (2° diameter, 2 cycles per degree,
668 2.5 Hz drift rate). Stimuli were generated with custom software (EXPO by P. Lennie) and displayed
669 on a cathode ray tube monitor (Hewlett Packard p1230; 1024×768 pixels, with ~ 40 cd/m² mean
670 luminance and 100 Hz frame rate) viewed at a distance of 110 cm (for anaesthetized dataset) or
671 60 cm (for awake dataset). Grating orientations were randomly interleaved, each presented for 70
672 ms (for anaesthetized dataset) or 150 ms (for awake dataset), separated by a uniform gray screen
673 (blank stimulus) for the same duration.

674 For each electrode, we extracted waveform signals (sampled at 30 kHz) whenever the extracel-
675 lular voltage exceeded a user defined threshold (typically 5x the root mean square signal on each
676 channel). We then sorted waveforms manually using the Plexon Offline Sorter, and isolated both
677 single and multi-unit clusters, here both referred to as neurons. We computed spike counts in a
678 fixed window with length equal to the stimulus duration, shifted by 50 ms after stimulus onset. We
679 excluded from further analyses all neurons that were not driven by any stimulus above baseline +
680 3std.

681 In the first dataset the monkey was awake, and there were $d_T = 3168$ trials of the responses of
682 72 neurons; due to the presence of cross-talk between a small subset of electrodes, we removed
683 all pairs of neurons in the dataset that exhibited correlations greater than 0.5, which left $d_N = 43$
684 neurons in the dataset. We refer to this dataset as the awake V1 dataset. After familiarization with
685 the restraining chair, headpost surgery, and postoperative recovery time (methods and protocols

686 described in *Festa et al. (2020)*), the animal was trained to fixate in a $1.3^\circ \times 1.3^\circ$ window. Eye posi-
687 tion was monitored with a high-speed infrared camera (Eyelink, 1000 Hz). A second surgery was
688 performed over V1 to implant a 96 channel microelectrode array into V1 (electrode length 1 mm).
689 After postoperative recovery, the spatial receptive fields of the sampled neurons were mapped by
690 presenting small patches of drifting full contrast gratings (0.5° diameter; 4 orientations, 1 cycle per
691 degree, 3 Hz drift rate, 250 ms presentation) at 25 distinct positions spanning a $3^\circ \times 4^\circ$ region of
692 visual space. Subsequent stimuli were centered in the aggregate receptive field of the recorded
693 units.

694 In the second dataset the monkey was anaesthetized and there were $d_T = 10,800$ trials of the
695 responses of $d_N = 70$ neurons; we refer to this dataset as the anaesthetized V1 dataset. The
696 protocol and general methods employed for the anaesthetized experiment have been described
697 previously (*Smith and Kohn, 2008*). In short, anaesthesia was induced with ketamine (10 mg/kg) and
698 maintained during surgery with isoflurane (1.5–2.5% in 95% O₂), switching to sufentanil (6–18 μ g/kg
699 per h, adjusted as needed) during recordings. Eye movements were reduced using vecuronium
700 bromide (0.15 mg/kg per h). Temperature was maintained in the 36–37 C° range, and relevant vital
701 signs (EEG, ECG, blood pressure, end-tidal PCO₂, temperature, and lung pressure) were monitored
702 continuously to ensure sufficient level of anaesthesia and well-being. A 10×10 multielectrode array
703 (400 μ m spacing, 1 mm length) was implanted into the upper layers of primary visual cortex, at a
704 depth of 0.6–0.8 mm.

705 All procedures were approved by the Albert Einstein College of Medicine and followed the guide-
706 lines in the United States Public Health Service Guide for the Care and Use of Laboratory Animals.

707 Acknowledgments

708 We would like to thank all the members of the labs of Ruben Coen-Cagli and Adam Kohn for their
709 regular feedback and support.

710 Competing interests

711 The authors declare they have no competing interests.

712 References

- 713 **Abbott LF**, Dayan P. The Effect of Correlated Variability on the Accuracy of a Population Code. *Neural compu-*
714 *tation*. 1999; 11(1):91–101.
- 715 **Amari Si**, Nagaoka H. *Methods of Information Geometry*, vol. 191. American Mathematical Soc.; 2007.
- 716 **Archer EW**, Koster U, Pillow JW, Macke JH. Low-Dimensional Models of Neural Population Activity in Sensory
717 Cortical Circuits. In: Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, editors. *Advances in*
718 *Neural Information Processing Systems 27* Curran Associates, Inc.; 2014.p. 343–351.
- 719 **Arieli A**, Sterkin A, Grinvald A, Aertsen A. Dynamics of Ongoing Activity: Explanation of the Large Variability in
720 Evoked Cortical Responses. *Science*. 1996 Sep; 273(5283):1868–1871.
- 721 **Bartolo R**, Saunders RC, Mitz AR, Averbek BB. Information-Limiting Correlations in Large Neural Populations.
722 *Journal of Neuroscience*. 2020 Feb; 40(8):1668–1678.
- 723 **Beck J**, Ma WJ, Latham PE, Pouget A. Probabilistic Population Codes and the Exponential Family of Distributions.
724 *Progress in Brain Research*. 2007; 165:509–519.
- 725 **Beck J**, Latham P, Pouget A. Marginalization in Neural Circuits with Divisive Normalization. *The Journal of*
726 *Neuroscience*. 2011; 31(43):15310–15319.
- 727 **Beck J**, Bejjanki VR, Pouget A. Insights from a Simple Expression for Linear Fisher Information in a Recurrently
728 Connected Population of Spiking Neurons. *Neural Computation*. 2011 Mar; 23(6):1484–1502.
- 729 **Bondy AG**, Haefner RM, Cumming BG. Feedback Determines the Structure of Correlated Variability in Primary
730 Visual Cortex. *Nature Neuroscience*. 2018 Apr; 21(4):598–606.

- 731 **Brunel N**, Nadal JP. Mutual Information, Fisher Information, and Population Coding. *Neural Computation*.
732 1998; 10(7):1731–1757.
- 733 **Chaniialidis C**, Evers L, Neocleous T, Nobile A. Efficient Bayesian Inference for COM-Poisson Regression Models.
734 *Statistics and Computing*. 2018 May; 28(3):595–608.
- 735 **Cover TM**, Thomas JA. *Elements of Information Theory*. 2nd ed ed. Hoboken, NJ: Wiley-Interscience; 2006.
- 736 **Dayan P**, Abbott LF. *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*.
737 Massachusetts Institute of Technology Press; 2005.
- 738 **Doya K**. *Bayesian Brain: Probabilistic Approaches to Neural Coding*. MIT press; 2007.
- 739 **Drugowitsch J**, Mendonça AG, Mainen ZF, Pouget A. Learning Optimal Decisions with Confidence. *Proceedings*
740 *of the National Academy of Sciences*. 2019 Dec; 116(49):24872–24880.
- 741 **Ecker AS**, Berens P, Cotton RJ, Subramaniyan M, Denfield GH, Cadwell CR, Smirnakis SM, Bethge M, Tolias AS.
742 State Dependence of Noise Correlations in Macaque Primary Visual Cortex. *Neuron*. 2014 Apr; 82(1):235–
743 248.
- 744 **Ecker AS**, Berens P, Tolias AS, Bethge M. The Effect of Noise Correlations in Populations of Diversely Tuned
745 Neurons. *Journal of Neuroscience*. 2011 Oct; 31(40):14272–14283.
- 746 **Ecker AS**, Denfield GH, Bethge M, Tolias AS. On the Structure of Neuronal Population Activity under Fluctuations
747 in Attentional State. *The Journal of Neuroscience*. 2016 Feb; 36(5):1775–1789.
- 748 **Festa D**, Aschner A, Davila A, Kohn A, Coen-Cagli R. Neuronal Variability Reflects Probabilistic Inference Tuned
749 to Natural Image Statistics. *bioRxiv*. 2020 Jun; p. 2020.06.17.142182.
- 750 **Ganguli D**, Simoncelli EP. Efficient Sensory Encoding and Bayesian Inference with Heterogeneous Neural Pop-
751 ulations. *Neural Computation*. 2014 Oct; 26(10):2103–2134.
- 752 **Ganmor E**, Segev R, Schneidman E. A Thesaurus for a Neural Population Code. *eLife*. 2015 Sep; 4:e06134.
- 753 **Goris RL**T, Movshon JA, Simoncelli EP. Partitioning Neuronal Variability. *Nature Neuroscience*. 2014 Jun;
754 17(6):858–865.
- 755 **Graf ABA**, Kohn A, Jazayeri M, Movshon JA. Decoding the Activity of Neuronal Populations in Macaque Primary
756 Visual Cortex. *Nature Neuroscience*. 2011 Feb; 14(2):239–245.
- 757 **Granot-Atedgi E**, Tkačik G, Segev R, Schneidman E. Stimulus-Dependent Maximum Entropy Models of Neural
758 Population Codes. *PLOS Computational Biology*. 2013 Mar; 9(3):e1002922.
- 759 **Herz AV**, Mathis A, Stemmler M. Periodic Population Codes: From a Single Circular Variable to Higher Dimen-
760 sions, Multiple Nested Scales, and Conceptual Spaces. *Current Opinion in Neurobiology*. 2017 Oct; 46:99–
761 108.
- 762 **Inouye DI**, Yang E, Allen GI, Ravikumar P. A Review of Multivariate Distributions for Count Data Derived from the
763 Poisson Distribution: A Review of Multivariate Distributions for Count Data. *Wiley Interdisciplinary Reviews:*
764 *Computational Statistics*. 2017 May; 9(3):e1398.
- 765 **Kafashan M**, Jaffe A, Chettih SN, Nogueira R, Arandia-Romero I, Harvey CD, Moreno-Bote R, Drugowitsch J.
766 Scaling of Information in Large Neural Populations Reveals Signatures of Information-Limiting Correlations.
767 *bioRxiv*. 2020 Jan; p. 2020.01.10.902171.
- 768 **Kanashiro T**, Ocker GK, Cohen MR, Doiron B. Attentional Modulation of Neuronal Variability in Circuit Models
769 of Cortex. *eLife*. 2017 Jun; 6:e23978.
- 770 **Kanitscheider I**, Coen-Cagli R, Kohn A, Pouget A. Measuring Fisher Information Accurately in Correlated Neural
771 Populations. *PLoS computational biology*. 2015; 11(6):e1004218.
- 772 **Kanitscheider I**, Coen-Cagli R, Pouget A. Origin of Information-Limiting Noise Correlations. *Proceedings of the*
773 *National Academy of Sciences*. 2015 Dec; 112(50):E6973–E6982.
- 774 **Karlis D**, Meligkotsidou L. Finite Mixtures of Multivariate Poisson Distributions with Application. *Journal of*
775 *Statistical Planning and Inference*. 2007 Jun; 137(6):1942–1960.
- 776 **Kingma D**, Ba J. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*. 2014; .

- 777 **Kohn A**, Coen-Cagli R, Kanitscheider I, Pouget A. Correlations and Neuronal Population Information. Annual
778 Review of Neuroscience. 2016 Jul; 39(1):237–256.
- 779 **Kriegeskorte N**, Douglas PK. Cognitive Computational Neuroscience. Nature Neuroscience. 2018 Aug; p. 1.
- 780 **Lyamzin DR**, Macke JH, Lesica NA. Modeling Population Spike Trains with Specified Time-Varying Spike Rates,
781 Trial-to-Trial Variability, and Pairwise Signal and Noise Correlations. Frontiers in Computational Neuro-
782 science. 2010; 4.
- 783 **Ma WJ**, Beck J, Latham P, Pouget A. Bayesian Inference with Probabilistic Population Codes. Nature Neuro-
784 science. 2006 Oct; 9(11):1432–1438.
- 785 **Macke JH**, Buesing L, Cunningham JP, Yu BM, Shenoy KV, Sahani M. Empirical Models of Spiking in Neural
786 Populations. In: Shawe-Taylor J, Zemel RS, Bartlett PL, Pereira F, Weinberger KQ, editors. *Advances in Neural
787 Information Processing Systems 24* Curran Associates, Inc.; 2011.p. 1350–1358.
- 788 **Macke JH**, Murray I, Latham PE. How Biased Are Maximum Entropy Models? In: Shawe-Taylor J, Zemel RS,
789 Bartlett PL, Pereira F, Weinberger KQ, editors. *Advances in Neural Information Processing Systems 24* Curran
790 Associates, Inc.; 2011.p. 2034–2042.
- 791 **Makin JG**, Dichter BK, Sabes PN. Learning to Estimate Dynamical State with Probabilistic Population Codes.
792 PLoS Comput Biol. 2015 Nov; 11(11):e1004554.
- 793 **Maunsell JHR**. Neuronal Mechanisms of Visual Attention. Annual Review of Vision Science. 2015 Nov; 1(1):373–
794 391.
- 795 **McLachlan GJ**, Lee SX, Rathnayake SI. Finite Mixture Models. Annual Review of Statistics and Its Application.
796 2019; 6(1):355–378.
- 797 **Meshulam L**, Gauthier JL, Brody CD, Tank DW, Bialek W. Collective Behavior of Place and Non-Place Neurons
798 in the Hippocampal Network. Neuron. 2017 Dec; 96(5):1178–1191.e4.
- 799 **Montijn JS**, Liu RG, Aschner A, Kohn A, Latham PE, Pouget A. Strong Information-Limiting Correlations in Early
800 Visual Areas. bioRxiv. 2019 Nov; .
- 801 **Moreno-Bote R**, Beck J, Kanitscheider I, Pitkow X, Latham P, Pouget A. Information-Limiting Correlations. Na-
802 ture Neuroscience. 2014 Oct; 17(10):1410–1417.
- 803 **Neal RM**, Hinton GE. A View of the EM Algorithm That Justifies Incremental, Sparse, and Other Variants. In:
804 *Learning in Graphical Models* Springer; 1998.p. 355–368.
- 805 **Okun M**, Steinmetz NA, Cossell L, Iacaruso MF, Ko H, Barthó P, Moore T, Hofer SB, Mrsic-Flogel TD, Carandini M,
806 Harris KD. Diverse Coupling of Neurons to Populations in Sensory Cortex. Nature. 2015 May; 521(7553):511–
807 515.
- 808 **Panzeri S**, Harvey CD, Piasini E, Latham PE, Fellin T. Cracking the Neural Code for Sensory Perception by Com-
809 bining Statistics, Intervention, and Behavior. Neuron. 2017 Feb; 93(3):491–507.
- 810 **Park IM**, Meister MLR, Huk AC, Pillow JW. Encoding and Decoding in Parietal Cortex during Sensorimotor
811 Decision-Making. Nature Neuroscience. 2014 Oct; 17(10):1395–1403.
- 812 **Pillow JW**, Ahmadian Y, Paninski L. Model-Based Decoding, Information Estimation, and Change-Point Detec-
813 tion Techniques for Multineuron Spike Trains. Neural Computation. 2010 Oct; 23(1):1–45.
- 814 **Pillow JW**, Shlens J, Paninski L, Sher A, Litke AM, Chichilnisky EJ, Simoncelli EP. Spatio-Temporal Correlations
815 and Visual Signalling in a Complete Neuronal Population. Nature. 2008 Aug; 454(7207):995–999.
- 816 **Pitkow X**, Angelaki DE. Inference in the Brain: Statistics Flowing in Redundant Population Codes. Neuron. 2017
817 Jun; 94(5):943–953.
- 818 **Pouget A**, Drugowitsch J, Kepecs A. Confidence and Certainty: Distinct Probabilistic Quantities for Different
819 Goals. Nature neuroscience. 2016; 19(3):366–374.
- 820 **Rabinowitz NC**, Goris RL, Cohen M, Simoncelli EP. Attention Stabilizes the Shared Gain of V4 Populations. eLife.
821 2015 Nov; 4:e08998.
- 822 **Ruff DA**, Cohen MR. Simultaneous Multi-Area Recordings Suggest That Attention Improves Performance by
823 Reshaping Stimulus Representations. Nature Neuroscience. 2019 Sep; p. 1–8.

- 824 **Rumyantsev OI**, Lecoq JA, Hernandez O, Zhang Y, Savall J, Chrapkiewicz R, Li J, Zeng H, Ganguli S, Schnitzer MJ.
825 Fundamental Bounds on the Fidelity of Sensory Cortical Coding. *Nature*. 2020 Mar; p. 1–6.
- 826 **Runyan CA**, Piasini E, Panzeri S, Harvey CD. Distinct Timescales of Population Coding across Cortex. *Nature*.
827 2017 Aug; 548(7665):92–96.
- 828 **Schneidman E**. Towards the Design Principles of Neural Population Codes. *Current Opinion in Neurobiology*.
829 2016 Apr; 37:133–140.
- 830 **Schneidman E**, Berry MJ, Segev R, Bialek W. Weak Pairwise Correlations Imply Strongly Correlated Network
831 States in a Neural Population. *Nature*. 2006 Apr; 440(7087):1007–1012.
- 832 **Semedo JD**, Zandvakili A, Machens CK, Yu BM, Kohn A. Cortical Areas Interact through a Communication Sub-
833 space. *Neuron*. 2019 Apr; 102(1):249–259.e4.
- 834 **Seriès P**, Latham PE, Pouget A. Tuning Curve Sharpening for Orientation Selectivity: Coding Efficiency and the
835 Impact of Correlations. *Nature neuroscience*. 2004; 7(10):1129.
- 836 **Shidara M**, Mizuhiki T, Richmond BJ. Neuronal Firing in Anterior Cingulate Neurons Changes Modes across
837 Trials in Single States of Multitrial Reward Schedules. *Experimental Brain Research*. 2005 May; 163(2):242–
838 245.
- 839 **Shmueli G**, Minka TP, Kadane JB, Borle S, Boatwright P. A Useful Distribution for Fitting Discrete Data: Re-
840 vival of the Conway–Maxwell–Poisson Distribution. *Journal of the Royal Statistical Society: Series C (Applied*
841 *Statistics)*. 2005; 54(1):127–142.
- 842 **Smith MA**, Kohn A. Spatial and Temporal Scales of Neuronal Correlation in Primary Visual Cortex. *Journal of*
843 *Neuroscience*. 2008 Nov; 28(48):12591–12603.
- 844 **Snow M**, Coen-Cagli R, Schwartz O. Specificity and Timescales of Cortical Adaptation as Inferences about Nat-
845 ural Movie Statistics. *Journal of Vision*. 2016 Oct; 16(13).
- 846 **Snow M**, Coen-Cagli R, Schwartz O. Adaptation in the Visual Cortex: A Case for Probing Neuronal Populations
847 with Natural Stimuli. *F1000Research*. 2017 Jul; 6:1246.
- 848 **Snyder AC**, Morais MJ, Kohn A, Smith MA. Correlations in V1 Are Reduced by Stimulation Outside the Receptive
849 Field. *Journal of Neuroscience*. 2014 Aug; 34(34):11222–11227.
- 850 **Sokoloski S**. Implementing a Bayes Filter in a Neural Circuit: The Case of Unknown Stimulus Dynamics. *Neural*
851 *Computation*. 2017 Jun; 29(9):2450–2490.
- 852 **Sokoloski S**. Implementing Bayesian Inference with Neural Networks. Dissertation, University of Leipzig; 2019.
- 853 **Sompolinsky H**, Yoon H, Kang K, Shamir M. Population Coding in Neuronal Systems with Correlated Noise.
854 *Physical Review E*. 2001 Oct; 64(5).
- 855 **Stevenson IH**. Flexible Models for Spike Count Data with Both Over- and under- Dispersion. *Journal of Com-*
856 *putational Neuroscience*. 2016 Aug; 41(1):29–43.
- 857 **Sur P**, Shmueli G, Bose S, Dubey P. Modeling Bimodal Discrete Data Using Conway-Maxwell-Poisson Mixture
858 Models. *Journal of Business & Economic Statistics*. 2015 Jul; 33(3):352–365.
- 859 **Taouali W**, Benvenuti G, Wallisch P, Chavane F, Perrinet LU. Testing the Odds of Inherent vs. Observed Overdis-
860 persion in Neural Spike Counts. *Journal of Neurophysiology*. 2015 Oct; 115(1):434–444.
- 861 **Wainwright MJ**, Jordan MI. Graphical Models, Exponential Families, and Variational Inference. *Foundations*
862 *and Trends® in Machine Learning*. 2008; 1(1-2):1–305.
- 863 **Walker EY**, Cotton RJ, Ma WJ, Tolias AS. A Neural Basis of Probabilistic Computation in Visual Cortex. *Nature*
864 *Neuroscience*. 2020 Jan; 23(1):122–129.
- 865 **Wei XX**, Stocker AA. A Bayesian Observer Model Constrained by Efficient Coding Can Explain ‘anti-Bayesian’
866 Percepts. *Nature Neuroscience*. 2015 Sep; 18(10):1509–1517.
- 867 **Whiteway MR**, Butts DA. The Quest for Interpretable Models of Neural Population Activity. *Current Opinion*
868 *in Neurobiology*. 2019 Oct; 58:86–93.

- 869 **Wiener MC**, Richmond BJ. Decoding Spike Trains Instant by Instant Using Order Statistics and the Mixture-of-
870 Poissons Model. *Journal of Neuroscience*. 2003 Mar; 23(6):2394–2406.
- 871 **Yamins DLK**, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-Optimized Hierarchical Models
872 Predict Neural Responses in Higher Visual Cortex. *Proceedings of the National Academy of Sciences*. 2014
873 Oct; 111(23):8619–8624.
- 874 **Yerxa TE**, Kee E, DeWeese MR, Cooper EA. Efficient Sensory Coding of Multidimensional Stimuli. *PLOS Compu-*
875 *tational Biology*. 2020 Sep; 16(9):e1008146.
- 876 **Yu BM**, Cunningham JP, Santhanam G, Ryu SI, Shenoy KV, Sahani M. Gaussian-Process Factor Analysis for
877 Low-Dimensional Single-Trial Analysis of Neural Population Activity. *Journal of Neurophysiology*. 2009 Jul;
878 102(1):614–635.
- 879 **Zemel RS**, Dayan P, Pouget A. Probabilistic Interpretation of Population Codes. *Neural computation*. 1998;
880 10(2):403–430.
- 881 **Zhao Y**, Park IM. Variational Latent Gaussian Process for Recovering Single-Trial Dynamics from Population
882 Spike Trains. *Neural Computation*. 2017 Mar; 29(5):1293–1316.
- 883 **Zohary E**, Shadlen MN, Newsome WT. Correlated Neuronal Discharge Rate and Its Implications for Psychophys-
884 ical Performance. *Nature*. 1994 Jul; 370(6485):140–143.