
ORIGINAL ARTICLE**Analysis**

Sensitivity and specificity of a Bayesian single trial analysis for time varying neural signals

Jeff T. Mohl¹⁻³ | Valeria C. Caruso¹⁻⁵ | Surya T. Tokdar^{1,6} | Jennifer M. Groh¹⁻⁴

¹Duke Institute for Brain Sciences, Duke University, Durham, NC, 27708, USA

²Center for Cognitive Neuroscience, Duke University

³Department of Neurobiology, Duke University

⁴Department of Psychology and Neuroscience, Duke University

⁵Center for Human Growth and Development, University of Michigan, Ann Arbor, MI, 48109, USA

⁶Department of Statistical Science, Duke University

Correspondence

Jeff T. Mohl, Department of Neurobiology, Duke University, Durham, NC, 27708, USA
Email: jeffrey.mohl@duke.edu

Funding information

This work was supported by an NDSEG research fellowship to JTM (32 CFR 168a) and by a grant from the NIH to JMG (R01 DC016363)

We recently reported the existence of fluctuations in neural signals that may permit neurons to code multiple simultaneous stimuli sequentially across time [1]. This required deploying a novel statistical approach to permit investigation of neural activity at the scale of individual trials. Here we present tests using synthetic data to assess the sensitivity and specificity of this analysis. We fabricated datasets to match each of several potential response patterns derived from single-stimulus response distributions. In particular, we simulated dual stimulus trial spike counts that reflected fluctuating mixtures of the single stimulus spike counts, stable intermediate averages, single stimulus winner-take-all, or response distributions that were outside the range defined by the single stimulus responses (such as summation or suppression). We then assessed how well the analysis recovered the correct response pattern as a function of the number of simulated trials and the difference between the simulated responses to each “stimulus” alone. We found excellent recovery of the mixture, intermediate, and outside categories (>97% correct), and good recovery of the single/winner-take-all category (>90% correct) when the

number of trials was >20 and the single-stimulus response rates were 50Hz and 20Hz respectively. Both larger numbers of trials and greater separation between the single stimulus firing rates improved categorization accuracy. These results provide a benchmark, and guidelines for data collection, for use of this method to investigate coding of multiple items at the individual-trial time scale.

KEYWORDS

statistical models, single trial analysis, validation, multiple stimuli

1 | INTRODUCTION

We recently showed that when multiple stimuli are present, some neurons exhibit activity patterns that fluctuate between those evoked by each stimulus alone [1]. This dynamic code could allow the representation of all stimuli within the same population of neurons. Such fluctuations may be a widespread phenomenon in the brain, but would be overlooked using conventional analysis methods that investigate mean activity pooled across trials. Of particular interest are cases in which the time-and-trial-pooled responses evoked by multiple stimuli appear to reflect the average of the responses to each stimulus presented in isolation. This phenomenon, known as divisive normalization [2], has been observed in visual brain areas such as V1 and MT [3] as well as other sensory and cognitive domains [4, 5, 6, 7, 8, 9]. However, such responses could either reflect a true averaging of the two stimuli/conditions, producing a consistent stable intermediate level of firing on each trial, or could reflect a dynamic code that flexibly shifts between the individual stimuli across trials.

To evaluate neural responses on a single trial basis, the novel statistical approach introduced in Caruso et al. (2018) characterizes the distribution of spike counts elicited in response to two simultaneous stimuli using Bayesian inference. Here we provide a general assessment of the sensitivity and specificity of that approach by simulating known neural responses as benchmark cases. In particular, we investigate how the analysis performs as we parametrically varied the data sample size (number of trials), the mean firing rate of responses, and the difference between spike counts across conditions.

We demonstrate that our approach accurately categorizes synthetic neural data into expected categories. The robustness of the results depends heavily on sample size, as well as on firing rate differences between the two single cue conditions. Importantly, the model performs very well under reasonable experimental values (20 trials per condition, 60% firing rate change between conditions). Finally, we show that that the model gracefully handles datasets that do not exactly match any of the tested hypotheses. These results demonstrate the viability of the analysis method and provide constraints for interpretation of actual neural data.

2 | EXPERIMENTAL RATIONALE AND PROCEDURES

2.1 | Neural encoding patterns to be assessed

For simplicity, our approach focused on the case of two simultaneously presented stimuli (dual stimuli) but can be extended to a larger number of stimuli. We consider an experimental setup in which a neuron's response is recorded in

three interleaved conditions: in the presence of a single stimulus “A”, a single stimulus “B”, or both stimuli A and B (“AB”). We considered four possible response distributions to dual stimuli, in relation to the distributions observed when only one stimulus is present (Figure 1).

1. Neurons might respond to only one of the stimuli, and do so consistently (i.e. respond to the same one) across trials. One way this could occur would be if only one stimulus is located in a neuron’s receptive field, but it might also apply when both stimuli are in the receptive field (sometimes referred to as a winner-take-all encoding). We label this possibility “single”.
2. The responses to dual stimuli might be greater than the maximum or less than the minimum of the single-stimulus responses. This category includes enhancement/summation, as well as suppression of the response to one stimulus by another. We refer to this case as “outside”.
3. The responses to dual stimuli are a consistent weighted average of the responses evoked by each stimulus alone. Here, the dual stimulus responses are between the bounds set by the two single stimulus responses, and cluster around a stable intermediate value. We refer to this case as “intermediate”.
4. The responses to dual stimuli may fluctuate such that on each trial the neuron appears to be responding to only one of the two stimuli. We term this possibility “mixture” because it reflects a mixture of two distributions of A and B stimulus responses. This is analogous to a winner-take-all except that the neuron is switching across trials rather than encoding the same stimulus each trial. Like the “intermediate” category, there could be a weighting factor such that a higher proportion of trials favor one stimulus over the other.

2.2 | Model construction, Bayesian model comparison, and synthetic data

These four possibilities can be formalized on the basis of how the spike distributions on combined stimulus trials AB compare to those observed when only A or B are presented alone. If A and B elicit spike counts according to Poisson distributions $Poi(\lambda_A)$ and $Poi(\lambda_B)$, then we can ask which of four competing hypotheses best describe the spike counts observed on combined AB trials:

- (a) Single: $F = Poi(\lambda)$ for either $\lambda = \lambda_A$ or $\lambda = \lambda_B$, with λ constant across trials
- (b) Outside: $F = Poi(\lambda)$ for some unknown $\lambda \notin [\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B)]$
- (c) Intermediate: $F = Poi(\lambda)$ for some unknown $\lambda \in [\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B)]$
- (d) Mixture: $F = \alpha \cdot Poi(\lambda_A) + (1 - \alpha) \cdot Poi(\lambda_B)$ for some unknown $\alpha \in (0, 1)$

The plausibility of each of these models was determined by computing the posterior probabilities of each model given the data, with a default Jeffreys’ prior [10] on each of the model specific rate (λ) parameters and on the mixing probability parameter (α). Each model was given a uniform prior probability (1/4) and posterior model probabilities were calculated by computation of relevant intrinsic Bayes factors [11] (see appendix A.1 for a thorough description of the models and model selection strategy).

To evaluate the sensitivity and specificity of this method, we built synthetic neuronal spiking datasets to match each of the four potential encoding strategies tested by the model. Consistent with our previous study [1], we focused on response patterns that could be modeled as deriving from Poisson distributions. In principle, the approach could be extended to other forms of response distributions, but this is beyond the scope of this work.

Data files were generated as spike times drawn using an independent Poisson point process sampled at 1 ms

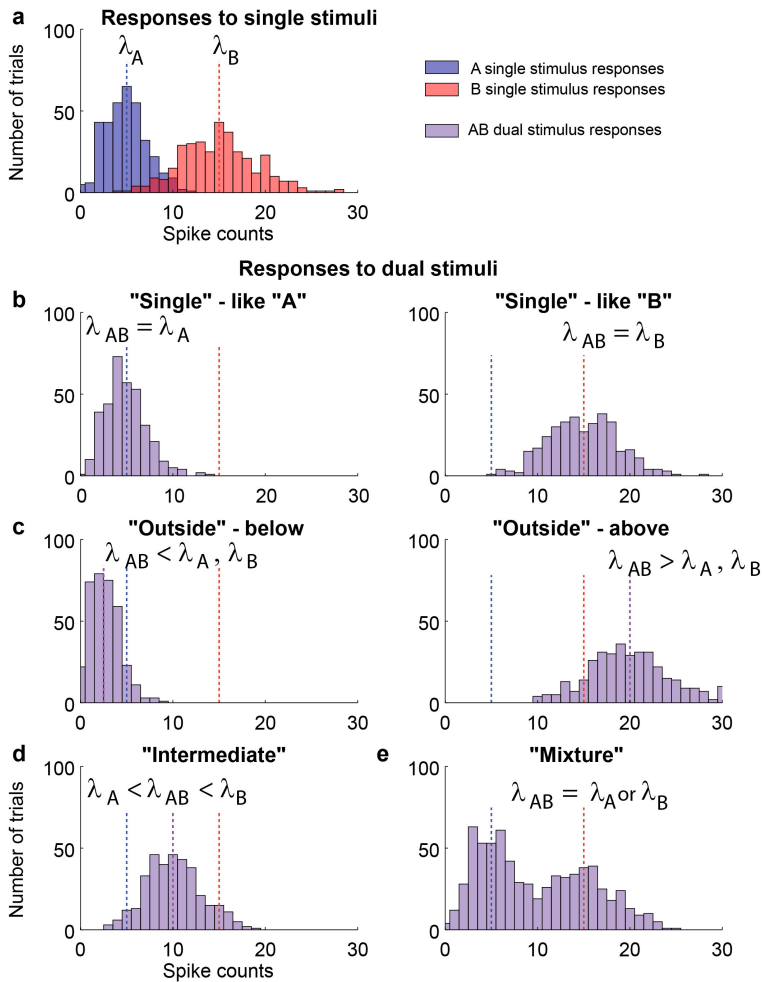


FIGURE 1 Four possible response patterns to dual stimuli trials, in relation to the responses observed to the component stimuli when presented individually. (a) Single stimulus trials were modeled as evoking spike counts distributed according to a Poisson process with rates λ_A , blue, or λ_B , red. (b) Responses on dual stimulus trials follow a Poisson λ_{AB} matching either λ_A , left, or λ_B , right. (c) The Poisson rate λ_{AB} on dual stimulus trials is less than (left) or greater than (right) those observed in single stimulus trials. In these simulations, λ_{AB} was set to $0.5 \cdot \lambda_A$ (left) or $\lambda_A + \lambda_B$ (right). (d) Spike counts derived from a Poisson process with a rate λ_{AB} between λ_A and λ_B . (e) Spike counts drawn from a mixture of two Poissons with rates λ_A and λ_B .

intervals, with constant mean firing rate for 1000 ms (Figure 2a-c). For A and B (single stimulus) trials, Poisson rates were assigned a priori to reflect a range of realistic firing rates for a single neuron presented with different stimuli. AB (combined stimulus) trials for each dataset were generated based on the chosen A and B firing rates and in a manner consistent with one of the four potential hypotheses. For the "single" hypothesis the AB data were generated using a single Poisson with rate λ_{AB} equal to the highest of the component rates, i.e. $\max(\lambda_A, \lambda_B)$. For the "outside" hypothesis, the rate λ_{AB} was set 20% higher than $\max(\lambda_A, \lambda_B)$. For the "intermediate" hypothesis, λ_{AB} was equal to the mean of

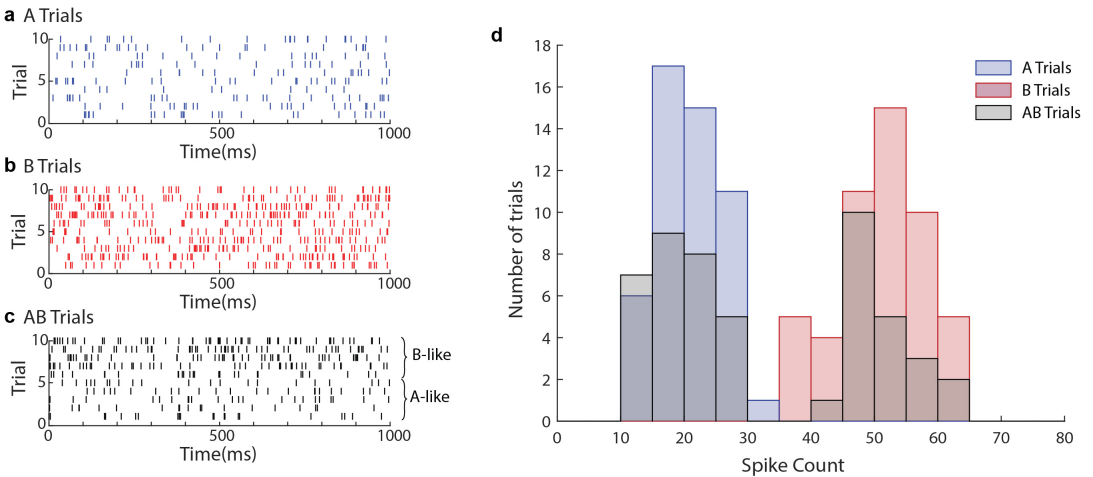


FIGURE 2 One example synthetic dataset. (a-c) Raster plots for a synthetic dataset built to match the across trial switching hypothesis, where blue rows (a) are single A trials, red rows (b) are single B trials, and black rows (c) are AB trials. AB trials are drawn randomly to match either A or B rates and are sorted so that B-like rates are towards the top of the raster. Even with sorting, this pattern is challenging to see with the naked eye, highlighting the need for analytical methods. (d) Whole trial spike count histogram for 50 repetitions of A, B, and AB trials. From this plot, the bimodality of AB trials for the switching condition is more apparent.

A and B rates $\lambda_{AB} = 0.5(\lambda_A) + 0.5(\lambda_B)$. For the across trial switching "mixture" model, the data were generated using the same Poisson process, but each trial was randomly chosen to be drawn from either $Poi(\lambda_A)$ or $Poi(\lambda_B)$ with equal probability. This results in a dataset for which the across trial average firing rate is equal to the average of the λ_A and λ_B rates, but individually each trial is better described as deriving from either the λ_A or λ_B response distributions. Note that it is nearly impossible to tell by visual inspection of a raster plot when a neuron has such a mixed response pattern, even when the trials are sorted as they are in Figure 2c, but the pattern becomes more evident in histograms of the trial-wise spike counts (Figure 2d).

Multiple datasets were generated using this strategy in order to test the power and reliability of the analysis under plausible experimental conditions. These datasets varied both the number of trials per condition (5-50 trials per condition) and the firing rates of the A and B conditions (from 5-100 Hz, with relative separation between A and B rates of 33-80% of maximum rate). Individual triplet pairs were generated under each of these conditions, analogous to running 100 individual cells through the analysis. This set of parameters was used for all conditions tested, including datasets constructed to not exactly match any of the hypotheses, discussed in the final section of the results.

2.3 | Code and data availability

Code specific to this paper can be found on GitHub at https://github.com/jmohl/mplx_tests, (archived DOI: 10.5281/zenodo.3508536) which includes the code used to generate synthetic data for this manuscript as well as all code needed to perform the neural mixture analysis. The exact synthetic data files used to generate plots are available upon request. Source code and documentation for the Neural Mixture Model available at <https://github.com/tokdarstat/Neural-Multiplexing>.

3 | RESULTS

3.1 | Neural Mixture Model accurately characterizes synthetic data

The desired analysis outcome is for the output to match the input. That is, data explicitly generated to match the single hypothesis should be correctly labeled as “single”, data generated as “outside” should be labeled “outside”, etc. Figure 3 illustrates that this is largely the case. The series of simulations shown here involved 20 A trials simulated with $\lambda_A = 20\text{Hz}$, 20 B trials with $\lambda_B = 50\text{ Hz}$, and 20 AB trials generated according to the various methods specified above. The “mixture” and “intermediate” categories perform exceptionally well, with 100/100 “mixture” and 99/100 “intermediate” datasets labeled correctly with >95% confidence (dark black bars). This distinction is critical, as these two conditions would produce exactly the same mean rate when averaging across trials, making them indistinguishable using typical neural analysis strategies which average across trials in order to reduce noise. “Single” and “outside” datasets were also correctly labeled in the majority of cases (90/100 and 97/100, respectively), although these hypotheses are not the focus of our analysis as they can be differentiated more easily using simpler statistical methods.

Although the category “single” was correctly identified as the best model for the dataset simulated under the “single” hypothesis 90% of the time, the posterior probability or confidence level did not reach the 95% level observed for the other models. This is due to the narrow definition of this category: response rates on AB trials must be indistinguishable from those occurring on either the A or B trials. All other categories include a range of possibilities which admits this hypothesis as a boundary case (i.e. a weighted average with the weight for A set to 1). Therefore, these models are all competitive in explaining data that is generated to match the “single” case, which explains the low posterior probability of this model. For this reason, it is better to consider the “single” category as reflective of a null hypothesis, where there is no interaction at all between stimuli.

3.2 | Dependence on number of trials and difference between A and B responses

The accuracy of this characterization depended on both the amount of data and the difference between the response distributions on A and B trials. The dependence on the number of trials is best appreciated when considering similar A and B response distributions, such as $\lambda_A = 20$ vs $\lambda_B = 30\text{Hz}$ shown in Figure 4a, which depicts the average posterior probability value for the correct model as a function of the number of trials. Even with this modest separation between the A and B response patterns, increasing the number of trials per condition allowed the analysis to better characterize the underlying rates, and therefore better discriminate between competing hypotheses. “Single”, “Intermediate” and “Mixture” had average posterior probability values >0.3 for $n=5$ trials, but performance improves steadily to average posterior probability values of >0.75 for $n=50$ trials. When response distributions were moderately separated, $\lambda_A = 20$ vs $\lambda_B = 50\text{Hz}$ (the same separation used in Figure 3), performance rose more rapidly for all models except “single”. At $n=5$, posterior probability values range from 0.4 for “single” to 0.8 for “mixture”. At $n=30$, posterior probability values equaled approximately 1 for “mixture”, “intermediate” and “outside”. Further increasing the firing rate separation to $\lambda_A = 20$ vs $\lambda_B = 100\text{Hz}$ resulted in very high posterior probabilities of >0.95 even at $n=5$ for “mixture” and “intermediate”; this level was achieved for “outside” at $n=10$.

These figures give a rough sense of the sensitivity of our analysis, demonstrating that the analysis becomes more reliable as more trials per condition are added until reaching asymptote around 30 trials/condition for a 50 Hz vs 20 Hz comparison (Figure 4b). Similarly, increasing the difference in spike count between A and B conditions also improves specificity in the analysis, allowing for accurate characterization with as few as 5 trials (Figure 4c). Although these data were constructed under ideal conditions (the data perfectly matches one of the tested hypotheses), they can be used as

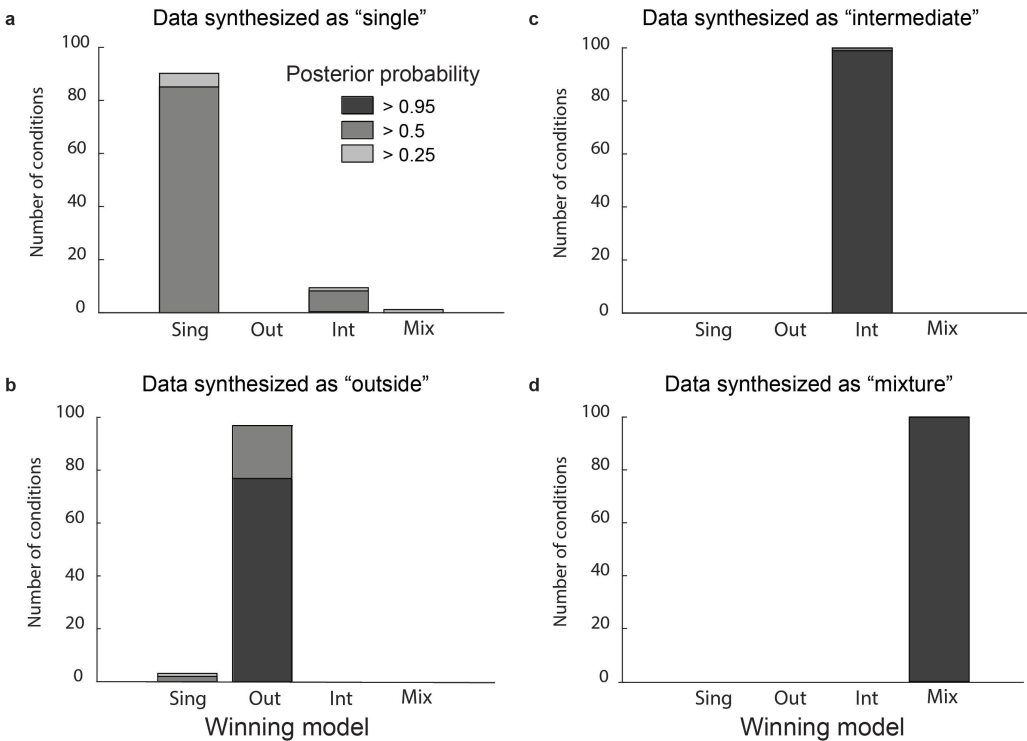


FIGURE 3 The analysis method correctly categorizes synthetic datasets created to match each model. The shading of the bars indicates the posterior probability with which each individual run of synthetic data ($n=100$) is assigned to a given category. Of particular interest is the very strong separation between intermediate and mixture datasets, as this discrimination is not possible when considering only firing rates averaged across trials. Parameters used for this figure: $\lambda_A = 20\text{Hz}$, $\lambda_B = 50\text{Hz}$, number of trials in each run = 20 per stimulus condition (60 overall).

a guide for how much data should be collected in order to obtain satisfactory results in a real dataset.

The above results give a detailed view of how each model performs across a realistic range of firing rates for neural recordings from primate sensory cortices and sub-cortical areas, for which we initially designed this method. We next sought to determine whether the analysis effectively extended into datasets with much lower maximum firing rates. To address this question we performed the analysis on data for a wide range of maximum firing rate values (5 to 100 Hz) for two fixed amounts of relative separation between A and B rates (40% and 80% of maximum rate) and characterized the prediction accuracy under each model (Figure 5). We found that firing rate affected the accuracy of the model, with lower average firing rate conditions resulting in worse performance than higher firing rates. As expected, a larger relative separation between A and B responses (analogous to having stronger neuronal preference for one or the other condition) resulted in significantly better performance, even for low firing rate conditions. However, even when considering the larger separation value of 80%, datasets with a maximum firing rate of $<15\text{Hz}$ barely reached 95% accuracy with 50 trials per condition. These results suggest that datasets with very low average firing rates (less than 15 Hz for the most preferred response) may not be resolvable using this analysis method.

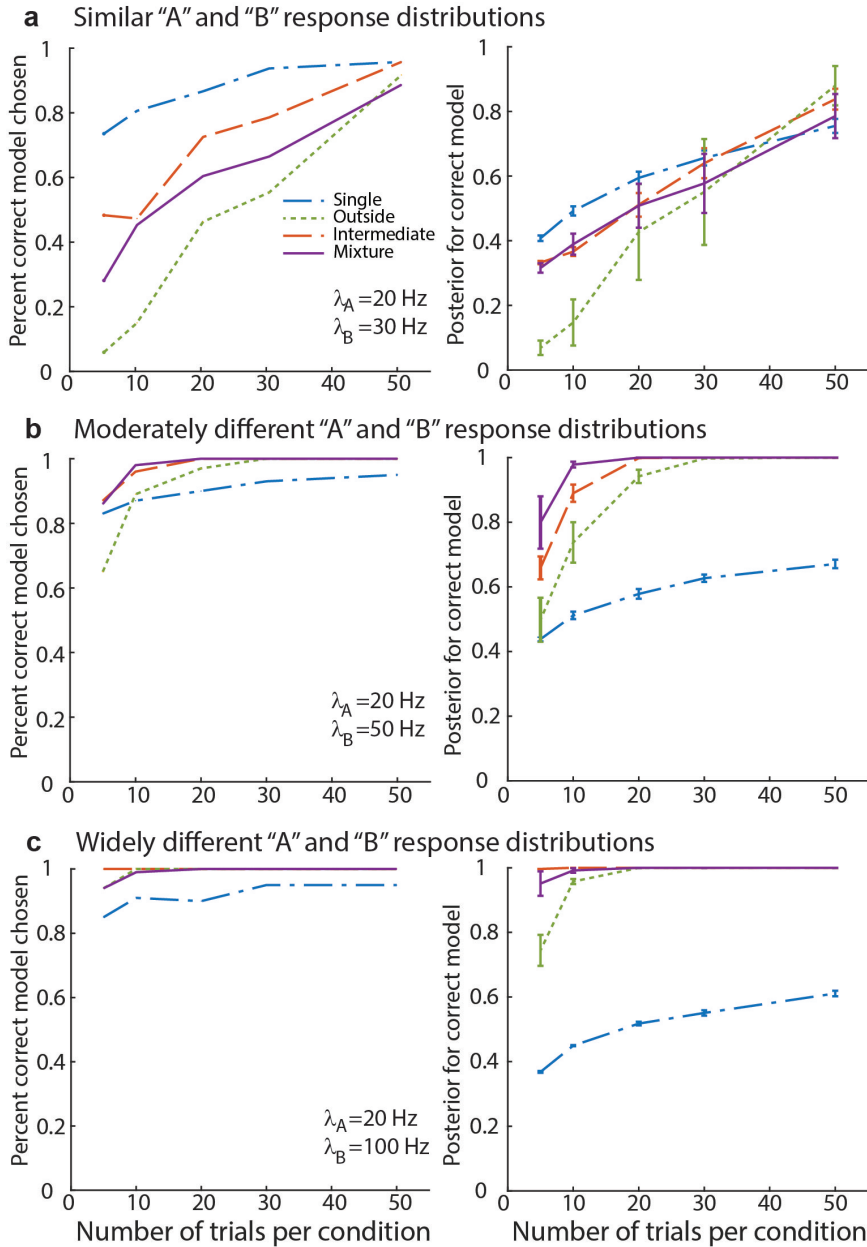


FIGURE 4 Increasing number of trials or separability of conditions improves accuracy of model comparison. (a) Left, percent of triplets which were correctly categorized, split by dataset type, for increasing number of trials per conditions; $\lambda_A = 20$ Hz, $\lambda_B = 30$ Hz. Right, mean and variance of posterior probability for correct model across triplets (b & c) same as in (a) but with λ_B set to 50 Hz and 100 Hz respectively. Fewer trials are needed when responses are very different between A and B trials.

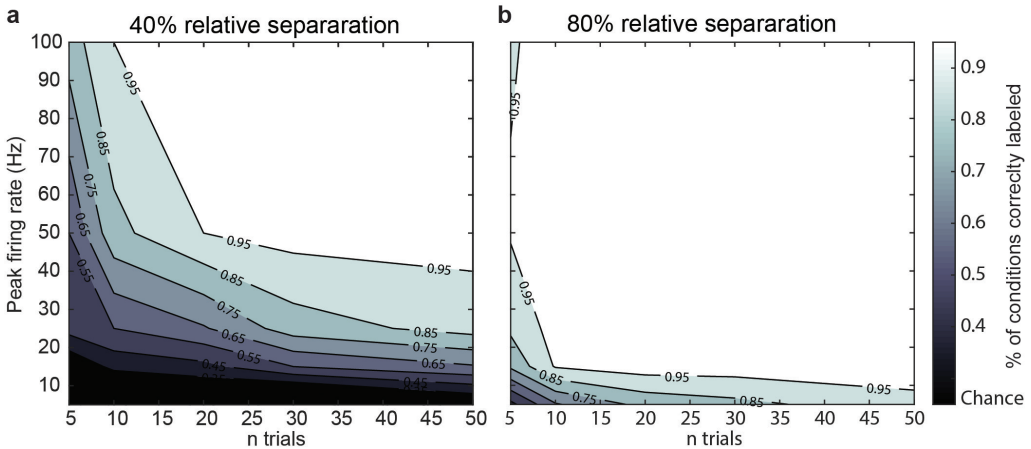


FIGURE 5 Model prediction accuracy depends on both number of trials and average firing rate. (a) Model classification accuracy collapsed across all four models for datasets generated for a fixed relative separation between A and B responses of 40% of the peak firing rate. Shading shows the percentage of conditions correctly categorized with interpolated phase transitions. (b) Same as in a, but for a relative separation of 80% of the peak rate.

3.3 | Model results are informative when datasets do not perfectly match hypotheses

Though this modeling strategy was meant to test between discrete hypotheses, it is unlikely that real neural signals perfectly and uniquely match any one of these scenarios. Therefore, it is important that the analysis accurately reflect deviations from exact hypothesis matches. Here, we consider two potential deviations from the circumstances considered above: weighted averaging of A and B stimuli and incomplete switching between A and B rates.

For weighted averaging datasets the AB trials were generated as in the “intermediate” condition above, except that the AB rate was set to be closer to the A rate than the B rate: $\lambda_{AB} = 0.75 \cdot \lambda_A + 0.25 \cdot \lambda_B$. Because the model returns both a classification and a posterior probability (reflecting the model’s confidence in that classification), we expected that this type of dataset will result in a spread across single and average classifications, but with lower confidence in this assessment. This is indeed the case, as the analysis returned primarily the intermediate category, with some single winners, but with much lower posterior probabilities than the well matched datasets (Figure 6a, compare with Figure 3c).

We also tested a form of incomplete switching, where stimuli show strong fluctuations but did not quite exactly match the A and B rates. These datasets were generated using the same strategy as the mixture datasets described above, except that A-like or B-like trials were generated with a slightly shifted mean firing rate. Multiple degrees of similarity were tested, but two (80% and 90% similarity) are presented here. From these, the analysis accurately described a 90% switch as mixture (Figure 6b), but around 80% similarity it began to interpret many datasets as average (Figure 6c). This highlights a natural limitation that should be expected in the data, as continuing to reduce the similarity would eventually result in a condition that was indistinguishable from true averaging. However, these results demonstrate the high specificity of our analysis for the mixture category, enforcing a strong definition of mixture (literally switching between rates closely matched to the A and B rates, rather than any amount of fluctuation).

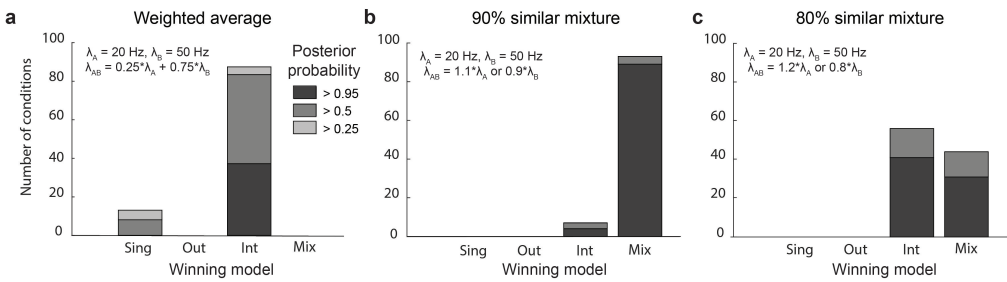


FIGURE 6 Datasets which do not exactly match the canonical hypotheses are descriptively categorized by model comparison. (a) A dataset generated to reflect weighted average of A and B stimuli is weakly categorized as intermediate with some triplets categorized as singles. (b) Mixture trials generated to alternate between values shifted 10% from the true A and B firing rates are primarily categorized as mixtures. (c) Mixture trials with rates shifted 20% from the true A and B rates are categorized as either mixture or average with equal probability, consistent with the fact that these trials would be much more difficult to discriminate from the true averaging hypothesis. Parameters used for this figure: $\lambda_B = 50\text{Hz}$, $\lambda_A = 20\text{Hz}$, number of trials=20.

4 | DISCUSSION

There is broad interest in understanding the nature and significance of firing patterns in the brain. It is well known that such firing patterns are variable in the face of identical, highly controlled experimental conditions (such as the presentation of the same stimulus in the same context). While many studies have viewed this variability as deleterious “noise” that if unsolved would undermine the ability of the brain to perform its essential tasks [12, 13, 14, 15, 16], we and others have sought explanations under the possibility that certain forms of such variation may contribute in a positive fashion to brain function [1, 17, 18, 19, 20, 21, 22, 23, 24]. In particular, we have successfully modeled variation in whole-trial spike counts to multiple stimuli as being drawn from the observed distributions of spike counts to those same stimuli when presented individually [1].

Here, we benchmarked this analysis on synthetic datasets to provide insight into the sensitivity and specificity of our analysis method as a function of trial counts and the separation between the distributions of spike counts elicited by the component individual stimuli. When response separations are large, e.g. the mean rate for one stimulus condition is 5X the rate for the other, the analysis method can successfully distinguish among the 4 competing hypotheses with as few as $n=5$ trials for each condition ($n=15$ overall). Smaller response separations can be compensated for by collecting more trials to achieve similarly good results. Finally, even when the conditions do not exactly match the assumptions, such as if the component response rates in the “mixture” condition do not precisely match those observed when the single stimuli were presented individually, correct classifications greatly outnumber incorrect ones. Critically, the analysis is conservative against the “mixture” hypothesis in these cases, demonstrating that data which is best fit by this model is truly fluctuating between the responses to single stimuli at the single trial level.

These results suggest that the analysis tested here are suitable for many electrophysiological datasets which match the A, B, AB format. Datasets which have moderately high peak firing rates (50 Hz) and an average response difference between conditions of approximately 40% (relative to peak rate) can reach over 95% categorization accuracy with as few as 20 trials. Higher peak firing rates, larger separation between neural response, or a greater number of trials all improve the accuracy of our analysis. Conversely, datasets with low peak firing rates (15 Hz) are likely to produce only weak results even with a large number of trials (which will be reflected in low model posterior probabilities). Practically, this means that our analysis is particularly well suited for recordings in primate sensory or motor brain regions with the

pronounced tuning and firing rate changes required to differentiate responses at the single trial level.

A situation not tested here is the case in which the response distributions are not derived from Poisson distributions. In our previous work [1] we excluded conditions in which the responses to individual stimuli did not satisfactorily resemble Poisson distributions in order to ensure that our model assumptions were not violated, but this has several downsides. First, it is difficult to have confidence in the success of this exclusion criterion: failing to reject the Poisson assumption is not the same as confirming its validity. Second, a considerable amount of data is excluded in this fashion (as much as 25-50%, depending on dataset, before even considering other exclusion criteria). Finally, there is significant evidence in the literature that spike counts in many brain areas are more variable than would be suggested by a Poisson distribution [14, 25, 26, 27, 28, 29]. For all of these reasons, it will be important to both test the model with data sets that violate this assumption and extend the analysis method to include other response distributions such as negative binomials.

The data presented here reflect conditions where two “stimuli” are presented at the same time, but this analysis could in principle be extended to combinations of three or more response patterns. We have previously shown that responses to multiple auditory stimuli in the primate inferior colliculus are often well described by mixtures of single stimulus responses [1], but little is known about how this type of code changes as more stimuli are added. An extension of this analysis into more complex mixtures of multiple different responses may help bring clarity to this question, and more work is needed to determine whether this phenomenon is general or limited to two stimulus cases.

Given the broad interest in both noise as a potential limitation on neural representations and in divisive normalization as an elemental computation in sensory processing – with recent suggestions that this process may be impaired in conditions such as autism [30, 31, 32] – it will be increasingly important to develop additional methods which can probe neural codes at the individual trial level [33, 34, 35]. The tools described in the present paper represent an important step towards uncovering fluctuating patterns in neural activity that may permit greater amounts of information to be encoded in the spike trains of individual and populations of neurons.

ACKNOWLEDGEMENTS

We would like to thank Shawn Willett and Meredith Schmell for help on an earlier version of this manuscript.

CONFLICT OF INTEREST

The authors have declared no competing interests related to this work.

REFERENCES

- [1] Caruso VC, Mohl JT, Glynn C, Lee J, Willett SM, Zaman A, et al. Single neurons may encode simultaneous stimuli by switching between activity patterns. *Nature Communications* 2018 dec;9(1):2715.
- [2] Carandini M, Heeger DJ. Summation and division by neurons in primate visual cortex. *Science* 1994;264(5163):1333–1336.
- [3] Britten KH, Heuer HW. Spatial summation in the receptive fields of MT neurons. *The Journal of neuroscience : the official journal of the Society for Neuroscience* 1999 jun;19(12):5074–84.
- [4] Bao P, Tsao DY. Representation of multiple objects in macaque category-selective areas. *Nature Communications* 2018;9(1).

- [5] Kozlov AS, Gentner TQ. Central auditory neurons have composite receptive fields. *Proceedings of the National Academy of Sciences* 2016;113(5):1441–1446.
- [6] Louie K, Khaw MW, Glimcher PW. Normalization is a general neural mechanism for context-dependent decision making. *Proceedings of the National Academy of Sciences* 2013;110(15):6139–6144.
- [7] Ohshiro T, Angelaki DE, DeAngelis GC. A normalization model of multisensory integration. *Nature neuroscience* 2011 jun;14(6):775–82.
- [8] Olsen SR, Bhandawat V, Wilson RI. Divisive normalization in olfactory population codes. *Neuron* 2010;66(2):287–299.
- [9] Reynolds JH, Heeger DJ. The Normalization Model of Attention. *Neuron* 2009;61(2):168–185.
- [10] Berger J. The case for objective Bayesian analysis. *Bayesian Analysis* 2006 sep;1(3):385–402.
- [11] Berger JO, Pericchi LR. The intrinsic bayes factor for model selection and prediction. *Journal of the American Statistical Association* 1996 mar;91(433):109–122.
- [12] Shadlen MN, Newsome WT. Noise, neural codes and cortical organization. *Current Opinion in Neurobiology* 1994;4(4):569–579.
- [13] Shadlen MN, Newsome WT. Is there a signal in the noise? *Current Opinion in Neurobiology* 1995;5(2):248–250.
- [14] Shadlen MN, Newsome WT. The variable discharge of cortical neurons: Implications for connectivity, computation, and information coding. *Journal of Neuroscience* 1998 may;18(10):3870–3896.
- [15] Softky WR, Koch C. The highly irregular firing of cortical cells is inconsistent with temporal integration of random EPSPs. *Journal of Neuroscience* 1993;13(1):334–350.
- [16] Zohary E, Shadlen MN, Newsome WT. Correlated neuronal discharge rate and its implications for psychophysical performance. *Nature* 1994;370(6485):140–143.
- [17] Akam T, Kullmann DM. Oscillatory multiplexing of population codes for selective communication in the mammalian brain. *Nature Reviews Neuroscience* 2014;15(2):111–122.
- [18] Hoppensteadt FC, Izhikevich EM. Thalamo-cortical interactions modeled by weakly connected oscillators: Could the brain use FM radio principles? *BioSystems* 1998;48(1-3):85–94.
- [19] Jezek K, Henriksen EJ, Treves A, Moser EI, Moser MB. Theta-paced flickering between place-cell maps in the hippocampus. *Nature* 2011;478(7368):246–249.
- [20] Li K, Kozyrev V, Kyllingsbæk S, Treue S, Ditlevsen S, Bundesen C. Neurons in Primate Visual Cortex Alternate between Responses to Multiple Stimuli in Their Receptive Field. *Frontiers in Computational Neuroscience* 2016 dec;10:141.
- [21] Lisman JE, Idiart MAP. Storage of 7 ± 2 short-term memories in oscillatory subcycles. *Science* 1995;267(5203):1512–1515.
- [22] Lisman JE, Jensen O. The Theta-Gamma Neural Code. *Neuron* 2013;77(6):1002–1016.
- [23] McLelland D, VanRullen R. Theta-Gamma Coding Meets Communication-through-Coherence: Neuronal Oscillatory Multiplexing Theories Reconciled. *PLoS Computational Biology* 2016;12(10).
- [24] Meister M. Multineuronal codes in retinal signaling. *Proceedings of the National Academy of Sciences* 2002;93(2):609–614.
- [25] Amarasingham A. Spike Count Reliability and the Poisson Hypothesis. *Journal of Neuroscience* 2006;26(3):801–809.

- [26] Barberini CL, Horwitz GD, Newsome WT. A Comparison of Spiking Statistics in Motion Sensing Neurones of Flies and Monkeys. In: *Motion Vision* Springer; 2011.p. 307–320.
- [27] Carandini M. Amplification of trial-to-trial response variability by neurons in visual cortex. *PLoS Biology* 2004;2(9).
- [28] Cur M, Beylin A, Snodderly DM. Response variability of neurons in primary visual cortex (V1) of alert monkeys. *Journal of Neuroscience* 1997 apr;17(8):2914–2920.
- [29] Recanzone GH. Rapidly induced auditory plasticity: the ventriloquism aftereffect. *Proceedings of the National Academy of Sciences of the United States of America* 1998 feb;95(3):869–875.
- [30] Rosenberg A, Patterson JS, Angelaki DE. A computational perspective on autism. *Proceedings of the National Academy of Sciences* 2015;112(30):9158–9165.
- [31] Rosenberg A, Sunkara A. Does attenuated divisive normalization affect gaze processing in autism spectrum disorder? A commentary on Palmer et al. (2018). *Cortex* 2019;111:316–318.
- [32] Van de Cruys S, Vanmarcke S, Steyaert J, Wagemans J. Intact perceptual bias in autism contradicts the decreased normalization model. *Scientific Reports* 2018;8(1).
- [33] Macke JH, Buesing L, Sahani M. Estimating state and parameters in state space models of spike trains. In: Chen Z, editor. *Advanced State Space Methods for Neural and Clinical Data* Cambridge: Cambridge University Press; 2015.p. 137–159.
- [34] Pandarinath C, O’Shea DJ, Collins J, Jozefowicz R, Stavisky SD, Kao JC, et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature Methods* 2018 oct;15(10):805–815.
- [35] Park IM, Meister MLRR, Huk AC, Pillow JW. Encoding and decoding in parietal cortex during sensorimotor decision-making. *Nat Neurosci* 2014 oct;17(10):1395–1403.
- [36] Berger JO, Guglielmi A. Bayesian and conditional frequentist testing of a parametric model versus nonparametric alternatives. *Journal of the American Statistical Association* 2001;96(453):174–184.
- [37] Tokdar ST, Kass RE. Importance sampling: a review. *Wiley Interdisciplinary Reviews: Computational Statistics* 2010;2(1):54–60.
- [38] Berger JO, Pericchi LR. The intrinsic Bayes factor for model selection and prediction. *Journal of the American Statistical Association* 1996;91(433):109–122.

A | SUPPLEMENTAL METHODS - MODEL DETAILS

A.1 | Introduction

Here we record in detail the modeling strategy used and specifics of the model selection procedure, the results of which are reported in the main text.

A.2 | Model

For each experimental condition $e \in \{A, B, AB\}$, let $Y_j^e, j = 1, \dots, n_e$ denote the spike counts from all n_e trials run under the condition. We model

1. $Y_j^A \sim \text{Poi}(\lambda_A), Y_j^B \sim \text{Poi}(\lambda_B)$ for some unknown $\lambda_A, \lambda_B > 0$; and,
2. $Y_j^{AB} \sim F$ with four competing hypotheses describing F
 - a. Mixture: $F = \alpha \cdot \text{Poi}(\lambda_A) + (1 - \alpha) \cdot \text{Poi}(\lambda_B)$ for some unknown $\alpha \in (0, 1)$
 - b. Intermediate: $F = \text{Poi}(\lambda)$ for some unknown $\lambda \in (\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B))$
 - c. Outside: $F = \text{Poi}(\lambda)$ for some unknown $\lambda \notin [\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B)]$
 - d. Single: $F = \text{Poi}(\lambda)$ for either $\lambda = \lambda_A$ or $\lambda = \lambda_B$, with the exact situation being unknown *a priori*.

A.3 | Method

A.3.1 | Bayesian testing

We carry out statistical testing between the set of hypotheses H listed above by adopting a Bayesian perspective. A prior probability p_h is assigned to each hypothesis $h \in H$, with $p_h > 0$ and $\sum_{h \in H} p_h = 1$. Let the observed data be denoted $Y = (Y_j^e : 1 \leq j \leq n_e, e \in \{A, B, AB\})$. Each hypothesis h gives rise to a model for Y which can be written generically as

$$Y \sim f_h(y | \theta_h), \theta_h \in \Theta_h$$

where θ_h captures all unknown parameters under the hypothesis. The modeling process is completed by assuming a prior distribution $\pi_h(\theta_h)$ on the parameter space Θ_h to reflect prior information and beliefs about the uncertainty about θ_h . Inference about θ_h is then drawn based on the uncertainty quantified by the resulting posterior distribution $\pi_h(\theta_h | Y) = \pi_h(\theta_h) f_h(Y | \theta_h) / f_h(Y)$ over Θ_h where the normalizing constant

$$f_h(Y) := \int_{\Theta_h} f_h(Y | \theta_h) \pi_h(\theta_h) d\theta_h$$

is recognized as the *marginal likelihood score* for hypothesis h given the observed data. Inference about the relative merits of the competing hypotheses is then drawn based on the posterior hypothesis probabilities

$$p_h(Y) = \frac{p_h f_h(Y)}{\sum_{h' \in H} p_{h'} f_{h'}(Y)}, h \in H, \quad (1)$$

which capture the post-data certainties about the competing hypotheses.

A.3.2 | Prior specification

Because the four competing hypotheses are only about the distribution of AB trial counts, and do not differ in their description of A and B trial count distributions, we adopt a common prior for the parameters pertaining to these latter distributions. Specifically, we take $\lambda_A \sim \text{Gam}(a, b)$ and $\lambda_B \sim \text{Gam}(a, b)$ for each of the four models.

For the Mixture hypothesis, the remaining model parameter is the mixing proportion $\alpha \in (0, 1)$. We assign it a beta prior: $\alpha \sim \text{Be}(c_1, c_2)$. For the Intermediate hypothesis, given λ_A and λ_B , the remaining parameter λ is assigned a conditional gamma prior $\text{Gam}(a, b)$ truncated to the interval $(\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B))$. Similarly, for the Outside hypothesis, we take the conditional prior on λ given λ_A, λ_B as the $\text{Gam}(a, b)$ distribution truncated to $(0, \infty) \cap [\min(\lambda_A, \lambda_B), \max(\lambda_A, \lambda_B)]^c$. In both these cases, the same a, b values are used as for the prior distributions for λ_A and λ_B .

A.3.3 | Computation

Computation of marginal likelihood scores $f_h(Y)$ is generally a complex task in Bayesian inference and require customized approaches to numerically evaluate the integration. Our prior choices and the low dimensionality of the parameter spaces associated with all the hypotheses make the task slightly easier for our problem. However, each hypothesis demands a different strategy to perform the integration and we give enough details below so that an enterprising student can implement these strategies from scratch.

Before getting into the details, we note a particular simplification that can be made to the expression of $p_h(Y)$ in (1) thanks to the assumption of a common prior distribution on λ_A and λ_B across all four hypotheses. Write $Y = (Y^A, Y^B, Y^{AB})$ where each Y^e denotes the data corresponding to experimental condition $e \in \{A, B, AB\}$. Then we can write

$$p_h(Y) = \frac{p_h \tilde{f}_h(Y^{AB} | Y^A, Y^B)}{\sum_{h' \in H} p_{h'} \tilde{f}_{h'}(Y^{AB} | Y^A, Y^B)} \quad (2)$$

where

$$\tilde{f}_h(Y^{AB} | Y^A, Y^B) = \int \left\{ \tilde{f}_h(Y^{AB} | \tilde{\theta}_h, \lambda_A, \lambda_B) \pi_h(\tilde{\theta}_h | \lambda_A, \lambda_B) \times \pi(\lambda_A | Y^A) \pi(\lambda_B | Y^B) \right\} d\tilde{\theta}_h d\lambda_A d\lambda_B$$

with \tilde{f}_h denoting the probability mass function of Y^{AB} under model h and $\tilde{\theta}_h$ denoting the remaining parameters of the model. Notice that

$$\pi(\lambda_A | Y^A) = \text{Gam}(a + S_A, b + n_A), \quad \pi(\lambda_B | Y^B) = \text{Gam}(a + S_B, b + n_B) \quad (3)$$

where $S_A = \sum_{j=1}^{n_A} Y_j^A$ and $S_B = \sum_{j=1}^{n_B} Y_j^B$.

A particular integration operation that shows up repeatedly in the following calculations stems from the well-know Poisson-Gamma conjugacy. For a vector of non-negative integers $y = (y_1, \dots, y_n)$ and positive real numbers α, β , we

define the quantity

$$g(y; \alpha, \beta) := \int \prod_{j=1}^n \text{Poi}(y_j | \lambda) \text{Gam}(\lambda | \alpha, \beta) d\lambda \quad (4)$$

$$= \frac{\Gamma(\alpha + S(y))}{\Gamma(\alpha)} \frac{\beta^\alpha}{(\beta + n)^{\alpha + S(y)}} \frac{1}{\prod_{j=1}^n y_j!} \quad (5)$$

where $S(y) = \sum_{j=1}^n y_j$. This quantity could be easily evaluated using any standard mathematics or statistics software. But we do want caution the user of numerical overflow problems as the gamma function grows super-exponentially in its argument. It is best to carry out the calculation of $g(y; \alpha, \beta)$ in the logarithmic scale, such as using the `lgamma()` function in the R software platform.

A.3.4 | Computation for “Single” hypothesis

Leveraging the Poisson-Gamma conjugacy, one can directly calculate

$$\begin{aligned} \tilde{f}(Y^{AB} | Y^A, Y^B) &= \frac{1}{2} \left[\int \prod_{j=1}^{n_{AB}} \text{Poi}(Y_j^{AB} | \lambda_A) \pi(\lambda_A | Y^A) d\lambda_A \right. \\ &\quad \left. + \int \prod_{j=1}^{n_{AB}} \text{Poi}(Y_j^{AB} | \lambda_B) \pi(\lambda_B | Y^B) d\lambda_B \right] \\ &= \frac{1}{2} \left[g(Y^{AB}; a + S_A, b + n_A) + g(Y^{AB}; a + S_B, b + n_B) \right]. \end{aligned}$$

The above calculation is done by assuming that the total prior probability of the Single hypothesis is split equally (*a priori*) between its two sub-hypotheses. A more conservative variation of this would be to report the maximum of the two numbers $g(Y^{AB}; a + S_A, b + n_A)$ and $g(Y^{AB}; a + S_B, b + n_B)$, giving full weight to the sub-hypothesis that explains the data better. Such selective representation of strongest sub-hypothesis has been used in the literature [36].

A.3.5 | Computation for “Mixture” hypothesis

For this hypothesis, the remaining parameter is $\tilde{\theta} = \alpha$ with $\pi(\alpha | \lambda_A, \lambda_B) = \text{Be}(c_1, c_2)$ and

$$\tilde{f}(Y^{AB} | \alpha, \lambda_A, \lambda_B) = \prod_{j=1}^{n_{AB}} \{ \alpha \cdot \text{Poi}(Y_j^{AB} | \lambda_A) + (1 - \alpha) \cdot \text{Poi}(Y_j^{AB} | \lambda_B) \}. \quad (6)$$

This form of \tilde{f} is difficult to work with directly because when the product of the sum is expanded, it results in too many summands; $2^{n_{AB}}$ many to be precise. Instead, a commonly adopted strategy in dealing with mixture model computation is to rewrite the model by introducing latent (unobserved) variables $Z_j \in \{A, B\}, j = 1, \dots, n_{AB}$ that indicate which of the two Poisson components observation j came from. By considering,

$$Z_j \sim \text{Discrete}(\{A, B\}; (\alpha, 1 - \alpha)); Y_j | (Z_j = c) \sim \text{Poi}(\lambda_c); j = 1, \dots, n_{AB}$$

we recover the same joint distribution \tilde{f} for Y^{AB} as in (6).

One may write $\tilde{f}(Y^{AB} | Y^A, Y^B) = \int \tilde{f}(Y^{AB} | Y^A, Y^B, Z) p(Z) dZ$ where $p(Z)$ denotes the joint distribution on

Z under the model (and the integral actually is a sum over a discrete space): $p(Z) = \int p(Z | \alpha)\pi(\alpha)d\alpha = B(c_1 + \#\{Z_j = A\}, c_2 + \#\{Z_j = B\})/B(c_1, c_2)$ where $B(\cdot, \cdot)$ is the beta function. The integral can be numerically approximated by importance sampling Monte Carlo as follows. Let $q(Z)$ denote any distribution on the space of Z . Then with $Z^m, m = 1, \dots, M$, denoting a large sample of independent draws of Z from $q(Z)$, one has

$$\tilde{f}(Y^{AB} | Y^A, Y^B) \approx \frac{1}{M} \sum_{m=1}^M \tilde{f}(Y^{AB} | Y^A, Y^B, Z = Z^m) \frac{p(Z = Z^m)}{q(Z = Z^m)}$$

by the strong law of large numbers. The quality of this Monte Carlo approximation is improved by choosing an *importance distribution* $q(Z)$ that closely resembles the posterior distribution $p(Z | Y^{AB}, Y^A, Y^B) \propto \tilde{f}(Y^{AB} | Y^A, Y^B, Z)p(Z)$; see [37] for more details. For our purposes, a good and convenient choice is a $q(Z)$ under which $Z_j \sim \text{Discrete}(\{A, B\}, (\bar{\alpha}_j, 1 - \bar{\alpha}_j)), j = 1, \dots, n$, where

$$\begin{aligned} \bar{\alpha}_j &= \frac{\int \text{Poi}(Y_j^{AB} | \lambda_A)\pi(\lambda_A | Y^A)d\lambda_A}{\int \text{Poi}(Y_j^{AB} | \lambda_A)\pi(\lambda_A | Y^A)d\lambda_A + \int \text{Poi}(Y_j^{AB} | \lambda_B)\pi(\lambda_B | Y^B)d\lambda_B}, \\ &= \frac{g(Y_j^{AB}; a + S_A, b + n_A)}{g(Y_j^{AB}; a + S_A, b + n_A) + g(Y_j^{AB}; a + S_B, b + n_B)}, \end{aligned}$$

which calculates the probability of classifying trial j as having come from condition A, under equal prior odds.

Therefore, to carry out the above importance sampling Monte Carlo, it is sufficient that we evaluate $\tilde{f}(Y^{AB} | Y^A, Y^B, Z)$. But this can be computed efficiently as

$$\begin{aligned} \tilde{f}(Y^{AB} | Y^A, Y^B, Z) &= \int \tilde{f}(Y^{AB} | \lambda^A, \lambda^B, Z)\pi(\lambda_A | Y^A)\pi(\lambda_B | Y^B)d\lambda_A d\lambda_B \\ &= \left\{ \int \prod_{j: Z_j=A} \text{Poi}(Y_j^{AB} | \lambda_A)\pi(\lambda_A | Y^A)d\lambda_A \right\} \times \left\{ \int \prod_{j: Z_j=B} \text{Poi}(Y_j^{AB} | \lambda_B)\pi(\lambda_B | Y^B)d\lambda_B \right\} \\ &= g(\{Y_j^{AB} : Z_j = A\}; a + S_A, b + n_A) \cdot g(\{Y_j^{AB} : Z_j = B\}; a + S_B, b + n_B) \end{aligned}$$

by using the Poisson-Gamma conjugacy.

A.3.6 | Computation for the “Intermediate” hypothesis

For the Intermediate hypothesis, one can use a straight Monte Carlo average to approximate $\tilde{f}(Y^{AB} | Y^A, Y^B)$ as

$$\begin{aligned} \tilde{f}(Y^{AB} | Y^A, Y^B) &= \int \left\{ \int \tilde{f}(Y^{AB} | \lambda)\pi(\lambda | \lambda_A, \lambda_B)d\lambda \right\} \pi(\lambda_A | Y^A)\pi(\lambda_B | Y^B)d\lambda_A d\lambda_B \\ &\approx \frac{1}{M} \sum_{m=1}^M \tilde{f}(Y^{AB} | \lambda_A = \lambda_A^m, \lambda_B = \lambda_B^m) \end{aligned}$$

where $(\lambda_A^m, \lambda_B^m)$, $m = 1, \dots, M$, are independent draws from $\pi(\lambda_A | Y^A) \times \pi(\lambda_B | Y^B)$ and, with $\underline{\lambda} = \min(\lambda_A, \lambda_B)$, $\bar{\lambda} = \max(\lambda_A, \lambda_B)$, $S_{AB} = \sum_{j=1}^{n_{AB}} Y_j^{AB}$,

$$\begin{aligned} \tilde{f}(Y^{AB} | \lambda_A, \lambda_B) &= \int \tilde{f}(Y^{AB} | \lambda) \pi(\lambda | \lambda_A, \lambda_B) d\lambda \\ &= \frac{\int_{\underline{\lambda}}^{\bar{\lambda}} \prod_{j=1}^n \text{Poi}(Y_j^{AB} | \lambda) \lambda^{a-1} e^{-b\lambda} d\lambda}{\int_{\underline{\lambda}}^{\bar{\lambda}} \lambda^{a-1} e^{-b\lambda} d\lambda} \\ &= \frac{\int_{\underline{\lambda}}^{\bar{\lambda}} \lambda^{a+S_{AB}-1} e^{-(b+n_{AB})\lambda} d\lambda}{\{\prod_{j=1}^n Y_j^{AB}!\} \int_{\underline{\lambda}}^{\bar{\lambda}} \lambda^{a-1} e^{-b\lambda} d\lambda} \\ &= g(Y^{AB}; a, b) \times \frac{F_{a+S_{AB}, b+n}(\bar{\lambda}) - F_{a+S_{AB}, b+n}(\underline{\lambda})}{F_{a,b}(\bar{\lambda}) - F_{a,b}(\underline{\lambda})} \end{aligned}$$

where $F_{\alpha, \beta}(x)$ is used to denote the cumulative distribution function of $\text{Gam}(\alpha, \beta)$.

A.3.7 | Computation for “Outside” hypothesis

Here the computation is done exactly as in the Intermediate hypothesis case, except for the following calculation which accounts for the fact that the conditional prior on λ given λ_A, λ_B is $\text{Gam}(a, b)$ truncated to the complement of the interval $(\underline{\lambda}, \bar{\lambda})$:

$$\tilde{f}(Y^{AB} | \lambda_A, \lambda_B) = g(Y^{AB}; a, b) \times \frac{1 - \{F_{a+S_{AB}, b+n}(\bar{\lambda}) - F_{a+S_{AB}, b+n}(\underline{\lambda})\}}{1 - \{F_{a,b}(\bar{\lambda}) - F_{a,b}(\underline{\lambda})\}}.$$

A.3.8 | Additional considerations for non-informative priors

An actual implementation of the above testing framework requires choosing the hyper-parameters a, b, c_1, c_2 , all positive valued real numbers. As with any Bayesian analysis, the results will have some dependence on the choice of these hyper-parameters. While expert knowledge about the model animal, brain region and sensory/cognitive task might help to choose reasonable values of these parameters, it may also be desirable to use some *default* values that encode minimal prior information about the model parameters.

One such approach is to use non-informative priors arising from Jeffreys' work. For the prior on the mixing proportion α , the Jeffreys prior is $Be(1/2, 1/2)$ which corresponds to our choice with $c_1 = c_2 = 1/2$. The Jeffreys' prior for estimating the mean of a Poisson distribution is the improper density function $\pi(\mu) \propto 1/\sqrt{\mu}, \mu > 0$, which matches, in a limiting sense, our choice of $\text{Gam}(a, b)$ with $a = 1/2$ and $b = 0$. This is because the posterior distribution for the Poisson mean under a $\text{Gam}(1/2, \beta)$ prior converges to the posterior distribution under the Jeffreys' prior as $\beta \rightarrow 0$.

However, such a limiting property does not hold for the marginal likelihood score! In fact, this score is not even well defined under the Jeffreys' prior, since prior density function is defined only up to a multiplicative constant. Also note that the quantity $g(y; \alpha, \beta) \rightarrow 0$ as $\beta \rightarrow 0$, and, hence working with a small but non-zero b is not an option either, since the resulting marginal likelihood scores for the Intermediate and the Outside hypotheses can be made arbitrarily small by choosing an arbitrarily small b .

Such anomalies can be effectively addressed by following the intrinsic Bayes factor approach of [38]. The Bayes factor between two hypotheses h and h' is defined as the ratio of the marginal likelihood scores $B_{h,h'}(Y) = f_h(Y)/f_{h'}(Y)$.

Notice that

$$\frac{p_h(Y)}{p_{h'}(Y)} = \frac{p_h}{p_{h'}} \times B_{h,h'}(Y), \quad (7)$$

that is the posterior odds between the two hypotheses depends on the data Y only through the Bayes factor. The intrinsic Bayes factor adjustment works for the case where data Y consists of a collection of observations (Y_1, \dots, Y_n) which, under each hypothesis h , are independently distributed with their distributions depending on a parameter $\theta_h \in \Theta_h$, with a prior distribution $\pi_h(\theta_h)$ chosen on Θ_h .

When one or both of $\pi_h(\theta_h)$ and $\pi_{h'}(\theta_{h'})$ are improper, defined only up to an arbitrary scaling factor, [38] recommend replacing them with proper distributions $\pi_h^\ell(\theta_h) = \pi_h(\theta_h | Y_\ell)$ and $\pi_{h'}^\ell(\theta_{h'}) = \pi_{h'}(\theta_{h'} | Y_\ell)$ obtained by calculating the posterior distribution given a small fraction of the data $Y_\ell = (Y_j : j \in \ell)$, for subset $\ell \subset \{1, \dots, n\}$ called a training set. A minimal training set is chosen, so that least amount of data is expended in this step of correcting for the impropriety of the prior distributions. Next, one calculates the marginal likelihood scores based on the new priors $\pi_h^\ell(\theta_h)$ and $\pi_{h'}^\ell(\theta_{h'})$, but using only the remaining part of the data $Y \setminus Y_\ell$. The resulting Bayes factor, which depends on the choice of the training set, but does not depend on any arbitrary scaling of the original priors, can be expressed as $B_{h,h'}^\ell(Y) = B_{h,h'}(Y)/B_{h,h'}(Y_\ell)$. To avoid the the effect of the arbitrary choice of the training set, one calculates the intrinsic Bayes factor $B_{h,h'}^*(Y)$ which is an average of $B_{h,h'}^\ell$ across all minimal training sets ℓ .

The final averaging could be an arithmetic, geometric or harmonic mean of the training set adjusted Bayes factors. We adopt the geometric mean approach, because it generalizes nicely to the case where one has more than two hypotheses to compare. The geometric mean intrinsic Bayes factor preserves the transitivity property that $B_{h,h'}^* = B_{h,h''}^* \times B_{h'',h'}^*$ and conforms to (7) with B replaced with B^* , for every pair of hypotheses $h, h' \in H$. Furthermore, the geometric mean intrinsic Bayes factor approach can be viewed as a direct adjustment to the marginal likelihood score $B_{h,h'}^*(Y) = f_h^*(Y)/f_{h'}^*(Y)$, where the corresponding intrinsic marginal likelihood score $f_h^*(Y)$ is defined as the geometric mean of $f_h(Y)/f_h(Y_\ell)$ across all minimal training sets ℓ .

For our four hypotheses, both Outside and Intermediate have improper priors when $b = 0$. For either hypothesis, a single observation is enough to give a proper posterior and hence the minimal training set size is one. Therefore the intrinsic marginal likelihood score adjustment for any or our models is achieved as:

$$\tilde{f}_h^*(Y^{AB} | Y^A, Y^B) = \frac{\tilde{f}_h(Y^{AB} | Y^A, Y^B)}{\left[\prod_{\ell=1}^{n_{AB}} \tilde{f}_h(Y_\ell^{AB} | Y^A, Y^B) \right]^{1/n_{AB}}} \quad (8)$$

where one uses the formulas derived above for $\tilde{f}_h(Y^{AB} | Y^A, Y^B)$ with $b \approx 0$. In our implementation we use $b = 10^{-5}$. The adjustments to the Single, Intermediate and Outside hypotheses are straightforward. For the Mixture model, one does not need to run an importance sampling Monte Carlo to calculate the denominator in (8). Instead, since for each $\ell \in \{1, \dots, n\}$ the corresponding Z_ℓ has only two possibilities $\{A, B\}$ a full enumeration can be done to express the denominator as $(c_1 + c_2)^{-1} \left[\prod_{\ell=1}^{n_{AB}} (c_1 g(Y_\ell^{AB}; a + S_A, b + n_A) + c_2 g(Y_\ell^{AB}; a + S_B, b + n_B)) \right]^{1/n_{AB}}$.