

---

# Toward Bayesian permutation inference for identifying neurons in *C. elegans*.

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

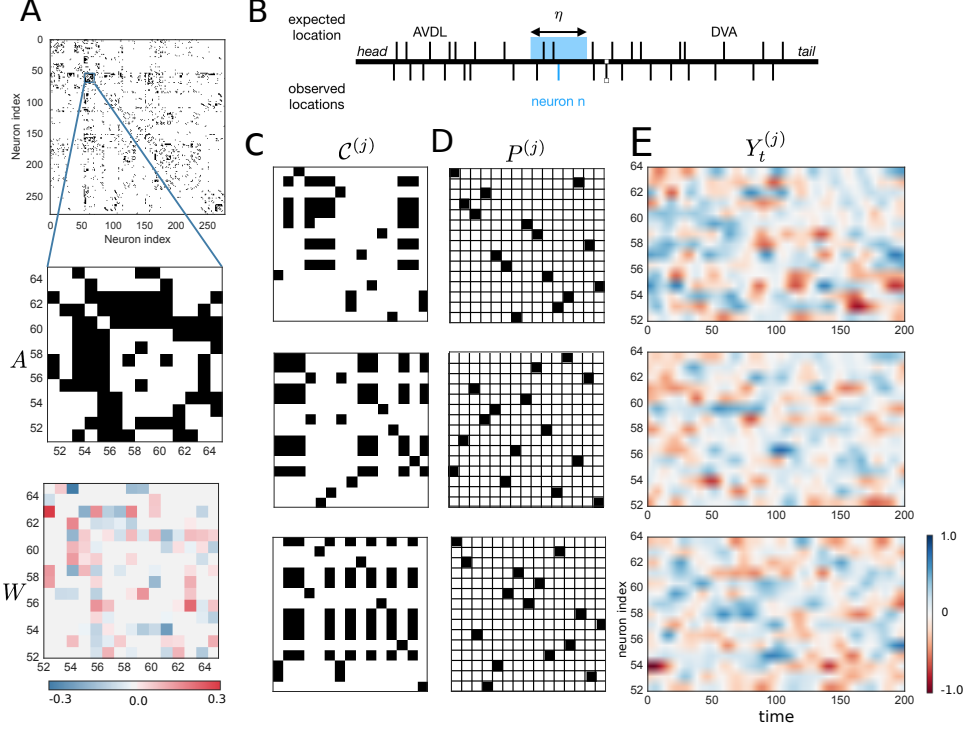
1 The nematode *C. elegans* is a unique model organism for neuroscientists as its  
2 connectome, or neural wiring diagram, has been known for at least three decades.  
3 Despite this knowledge, an understanding of the functional significance of these  
4 synaptic connections has remained elusive. Now several groups can routinely  
5 image the activity of a large fraction of neurons in the head of the worm, providing  
6 a unique opportunity to probe this organism. We propose a hierarchical Bayesian  
7 framework that combines strong prior information with data from many experi-  
8 ments to estimate posteriors over the functional connectivity weights. However,  
9 first we must clear a significant hurdle: in many cases we are not sure exactly  
10 which neurons are being imaged, so to combine information across experiments  
11 we must solve a matching, or permutation inference, problem. In this work we  
12 introduce new variational methods designed for joint inference of connectivity  
13 weights and neural identity. Working with actual permutations would involve  
14 evaluating and differentiating an intractable partition function. As an alternative,  
15 we build upon recent continuous relaxation techniques [Jang et al., 2016, Maddison  
16 et al., 2016], extending them from the original case of the probability simplex,  
17 to the Birkhoff polytope, the convex hull of permutation matrices. We test our  
18 method with simulated data from the true connectome and show that our approach  
19 outperforms many alternatives in identifying neurons.

## 20 1 Introduction

21 The nematode *C. elegans* plays a special role as a model organism in neuroscience since its neural  
22 network is stereotyped from animal to animal and its complete neural wiring diagram is known [Varsh-  
23 ney et al., 2011]. Modern calcium imaging technology enables measurements of hundreds of these  
24 neurons simultaneously [Kato et al., 2015, Nguyen et al., 2016]. The time is right to employ modern  
25 statistical methods to learn about the functional connectome in this system.

26 Ultimately, we are interested in the dynamical system that governs how neural activity evolves  
27 given its history and sensory inputs. Bayesian methods are ideally suited to this goal, allowing  
28 us to represent hierarchical probabilistic structures and integrate our prior knowledge about the  
29 connectome, the locations of neurons, etc. Bayesian learning and inference in dynamical systems  
30 with MCMC methods is well-studied, even for complicated models [De Freitas et al., 2001, Paninski  
31 et al., 2010]. Furthermore, hierarchical models to incorporate information from many worms are  
32 easily constructed in a Bayesian framework [Gelman et al., 2014].

33 However, our efforts to integrate information across worms are complicated by a major hurdle: in  
34 practice, associating recorded traces to neuron names is a painstaking, manual process. Experimenters  
35 consider the location of the neuron along with its pattern of activity to perform this matching, but the  
36 process is laborious and the results are prone to error. Without neuron names, we cannot represent



**Figure 1: Hierarchical Bayesian framework.** **A** We are given the actual adjacency matrix  $A$  from [Varshney et al., 2011]. The full matrix is shown (top) along with a zoom-in to 14 neurons (center). We wish to infer the corresponding weight matrix  $W$ , an example of which is shown below. **B** We also know the typical locations of the neurons [White et al., 1986, Lints et al., 2005]. Given observed locations, we constrain possible assignments to neuron identities within  $\eta$  of the observed location. **C** These constraints are represented as a matrix  $\mathcal{C}^{(j)}$  for worm  $j$  which specifies possible assignments of observed neurons to known identities. This illustration shows three worms. **D** To infer the weights, we must first infer the permutation  $P^{(j)}$  that matches the observed neurons in worm  $j$  to the set of known identities. **E** The observed data is a matrix  $Y_t^{(j)}$  whose rows are ordered according to the order in which neurons were observed in that worm. The permutation matrix maps this to the canonical ordering of the adjacency and weight matrices. Given  $\{Y_t^{(j)}\}_{j=1}^J$  and  $A$ , we infer  $\{P^{(j)}\}_{j=1}^J$  and  $W$ .

recordings canonically or learn about how one neuron influences another. This technical problem prevents the automatic use of hierarchical methods.

We present a method that aims to overcome this hurdle by incorporating inference over permutations that match observed neurons (*neuron 1*, *neuron 2*, ..., *neuron N*) to known names (*AVDL*, *AVAR*, ..., *SMDR*). Once the observed neurons have been mapped to canonical names, we can learn about the shared dynamical system. We focus on a simple linear autoregressive model for neural dynamics,

$$\tilde{Y}_t^{(j)} = (W \odot A) \tilde{Y}_{t-1}^{(j)} + \epsilon_t^{(j)}, \quad (1)$$

where  $W \in \mathbb{R}^{N \times N}$  is the weight matrix we wish to infer;  $A \in \{0, 1\}^{N \times N}$  is the known adjacency matrix or connectome;  $\odot$  denotes element-wise multiplication;  $\epsilon_t^{(j)} \sim \mathcal{N}(0, \sigma^2 I)$ ; and  $\tilde{Y}_t^{(j)} \in \mathbb{R}^N$  is the measured neural activity at time  $t$  in worm  $j$ . The catch is that  $\tilde{Y}_t^{(j)}$  is assumed to be in canonical order; i.e. in the same order as the rows and columns of  $W$  and  $A$ . We actually observe,

$$Y_t^{(j)} = P^{(j)} \tilde{Y}_t^{(j)}, \quad (2)$$

vectors that are permuted by matrix  $P^{(j)}$ . In order to learn about  $W$ , we must also infer the permutation matrices. We place a Gaussian prior on  $W$ .

The permutation matrices are constrained by side information. Specifically, we use neural position along the worm's body to constrain the possible neural identities for a given recorded neuron. We only allow an observed neuron to be mapped to a known identity if the observed location is within  $\eta$  of the

53 expected location. This is illustrated in Fig. 1B. We represent these constraints with the matrix  $\mathcal{C}^{(j)}$   
 54 so that  $\mathcal{C}_{mn}^{(j)} = 1$  if and only if observed neuron  $m$  is within  $\eta$  of canonical neuron  $n$ 's expected  
 55 location. An example is shown in Figure 1C. We let  $P^{(j)}$  have a uniform prior over the set of matrices  
 56 allowable under the given constraints.

57 We need to perform posterior inference of  $p(\{W, P^{(j)}\} | A, \{Y^{(j)}\})$ . MCMC with simple Metropolis-  
 58 Hastings proposals is straightforward, but we found this mixed poorly in practice. Motivated by recent  
 59 advances in automatic variational inference [Blei et al., 2017], we considered ways of extending this  
 60 technique to permutation inference. In Section 2 we detail our VI formulation and summarize the  
 61 methods we developed. Then in Section 3 we show that these methods outperform alternatives.

## 62 2 New methods for variational inference of latent permutations

63 Consider a latent variable model determined by a prior over the latent  $z \sim p(z)$  and a likelihood  
 64  $p(y | z)$  for the observed data  $y$ . In the VI framework, we approximate the intractable posterior  $p(z | y)$   
 65 with the distribution  $q \in \mathcal{Q}$  that best approximates the posterior. For tractability, we assume  $\mathcal{Q}$  is  
 66 indexed by a parameter  $\nu$ ; i.e.  $\mathcal{Q} = \{q(z; \nu) : \nu \in \mathcal{V}\}$ . The approximation is typically assessed  
 67 by the Kullback-Leibler (KL) divergence between the true posterior and variational approximation.  
 68 Minimizing the KL divergence is equivalent to maximizing the *evidence lower bound* (ELBO)

$$\mathcal{L}(\nu) \triangleq \mathbb{E}_{q(z; \nu)}[\log p(y|z)] - \text{KL}(q(z; \nu) \parallel p(z)), \quad (3)$$

69 with respect to  $\nu$ . We typically maximize equation (3) with stochastic optimization methods [Kushner  
 70 and Yin, 1987]: specifically, we approximate the expectations in (3) with Monte Carlo estimates  
 71 and optimize the ELBO with stochastic gradient ascent. One critical component is the choice of the  
 72 Monte Carlo approximation. Perhaps the most common choice is through the so called *score function*  
 73 *estimator*. Unfortunately, this estimator, also referred to as REINFORCE [Williams, 1992], cannot  
 74 be applied to permutations, since it involves the evaluation and differentiation of a likelihood which  
 75 is intractable for any non-trivial distribution over permutations (computing the partition function  
 76 involves a summation over  $N!$  terms).

77 The reparameterization trick Kingma and Welling [2013] offers an appealing alternative. If  $z$  can  
 78 be written as a differentiable function of a noise distribution and the parameters—i.e. if for certain  
 79  $f$  and  $\xi \sim p(\xi)$  one has  $z = f(\xi, \nu)$ —then we can write the expectation with respect to  $q(z)$  as an  
 80 expectation with respect to  $p(\xi)$  and bring the gradient inside the expectation. In the case of discrete  
 81 random variables a reparameterization always exists and it is given by the *Gumbel trick* [Papandreou  
 82 and Yuille, 2011], which states that one can sample from any discrete distribution by perturbing  
 83 each potential with Gumbel i.i.d noise, and then finding the configuration with the maximum value.  
 84 Unfortunately, the underlying  $f$  to this reparameterization is the non-differentiable  $\arg \max$  operator,  
 85 precluding the use of gradient descent methods.

86 Jang et al. [2016] and Maddison et al. [2016] proposed a solution to this problem, replacing the  
 87  $\arg \max$  by a temperature-dependent softmax approximation, which in the limit converges to the  
 88 original  $\arg \max$ . By combining the Gumbel trick with the softmax approximation, they conceived  
 89 the *Concrete* or *Gumbel-Softmax* distribution, and obtained explicit distribution formulae. Then,  
 90 they showed how to perform variational inference in discrete latent variable models using the  
 91 reparameterization trick and gradient descent. They replaced the original ELBO with a surrogate  
 92 appropriate for their continuous relaxation. The method works well if the temperature is chosen in a  
 93 reasonable range: not too high, to avoid degenerate distributions on the simplex; but also not too low,  
 94 to limit the variance the gradients.

95 We developed three methods for extending the Gumbel-softmax method to permutations. We name  
 96 then *stick-breaking*, *rounding* and *Gumbel-Sinkhorn* methods. We refer the reader to sections 3.1  
 97 and 3.2 of Linderman et al. [2017b] and section 4 of Anonymous [2018] for details, respectively.  
 98 Here we briefly summarize them: the primary geometric object is the Birkhoff polytope, the convex  
 99 hull of permutation matrices, and the analog of the probability simplex in the discrete case. In the  
 100 stick-breaking method, we generalize the standard construction on the simplex to stick-breaking of  
 101 the Birkhoff polytope. We show how to consistently “break the stick” while satisfying both the row  
 102 and column constraints that characterize a doubly stochastic matrix. For the rounding construction,  
 103 we start with a noise distribution and force it to be close to permutation matrices by rounding them  
 104 towards the extreme-points of the Birkhoff polytope (i.e. permutation matrices). Finally, for the

Table 1: Accuracy in the C.elegans neural identification problem, for varying mean number of candidate neurons (10, 30, 45, 60) and number of worms.

	10		30		45		60	
	1 worm	4 worms	1 Worm	4 worms	1 worm	4 worms	1 worms	4 worms
NAIVE VI	.34	.32	.16	.16	.13	.12	.11	.12
MAP	.34	.32	.17	.17	.14	.13	.13	.12
MCMC	.34	.65	.18	.28	.14	.17	.13	.15
VI	<b>.79</b>	<b>.94</b>	<b>.4</b>	<b>.69</b>	<b>.25</b>	<b>.51</b>	<b>.21</b>	<b>.44</b>

Table 2: Accuracy in inferring true neural identity for different of proportion of known neurons and  $\eta$ .

	40.%		30.%		20.%		10.%	
	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$
Naive VI	.43	.41	.33	.31	.23	.22	.12	.1
MAP	.42	.41	.33	.32	.23	.22	.12	.11
MCMC	.85	.80	.52	.46	.3	.26	.15	.12
VI	<b>.97</b>	<b>.96</b>	<b>.92</b>	<b>.84</b>	<b>.74</b>	<b>.58</b>	<b>.44</b>	<b>.23</b>

Gumbel-Sinkhorn method we notice that the so-called *Sinkhorn operator*, or infinite and successive row and column normalization of a matrix, is a natural extension of the softmax operator. With this, we introduce the Gumbel-Sinkhorn distribution, which approximates the sampling of a discrete distribution over permutations. Importantly, while stick-breaking and rounding yield explicit densities, Gumbel-Sinkhorn does not. However, there are ways to circumvent this difficulty, and overall we observe the latter performs the best.

### 3 Results

We evaluated various Bayesian inference methods for the hierarchical model illustrated in Figure 1. We compared against three alternatives: (i) naïve variational inference, where we do not enforce the constraint that  $P^{(j)}$  be a permutation and instead treat each row of  $P^{(j)}$  as a Dirichlet distributed vector; (ii) MCMC, where we alternate between sampling from the conditionals of  $W$  (Gaussian) and  $P^{(j)}$ , from which one can sample by proposing local swaps, as described in Diaconis [2009], and (iii) maximum a posteriori estimation (MAP). Our MAP algorithm alternates between the optimizing estimate of  $W$  given  $\{P^{(j)}, Y^{(j)}\}$  using linear regression and finding the optimal  $P^{(j)}$ . The second step requires solving a quadratic assignment problem (QAP) in  $P^{(j)}$ ; that is, it can be expressed as  $\text{Tr}(APBP^T)$  for matrices  $A, B$ . We used the QAP solver of Vogelstein et al. [2015].

We found that our method outperforms these alternative approaches. When there are many possible candidates (Table 1) and when only a small proportion of neurons are known with certitude (Table 2), variational inference via continuous relaxation with the Gumbel-Sinkhorn method performs best. Altogether, these results indicate our method enables a more efficient use of information than its alternatives. We conjecture that MCMC could be improved if local proposals—swapping pairs of labels—were replaced by more sophisticated transition operators, but fundamentally, it seems the hard assignments in the MCMC algorithm lead to poor mixing. We expect that the benefits of VI stem from the continuous relaxation, which enables soft assignments of neurons to identities.

Our results provide promising evidence that a Bayesian hierarchical approach to the study of neural dynamics on C. elegans is feasible. We note we made many simplifying assumptions that are not justified in practice: first, we assumed a linear dynamical system, while actual dynamics are highly nonlinear [Kato et al., 2015]. Fortunately, there exist many methods for inference in nonlinear system [Krishnan et al., 2015, Linderman et al., 2017a]. Also, we assumed all neurons were observed, while in reality we only see about 100 neurons at a time. Therefore, methods for marginalizing those unobserved set of neurons. The methods of Soudry et al. [2015] may help infer the weights, but reasoning about partial permutations requires more care. In conclusion, we have proposed a hierarchical Bayesian approach to the challenging neural identification problem, and while more work is needed, our initial results are promising.

## 139 References

- 140 Anonymous. Learning latent permutations with Gumbel-Sinkhorn networks. *International Confer-*  
141 *ence on Learning Representations*, 2018.
- 142 D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians.  
143 *Journal of the American Statistical Association*, 2017.
- 144 N. De Freitas, C. Andrieu, P. Højten-Sørensen, M. Niranjana, and A. Gee. Sequential Monte Carlo  
145 methods for neural networks. In *Sequential Monte Carlo methods in practice*, pages 359–379.  
146 Springer, 2001.
- 147 P. Diaconis. The Markov chain Monte Carlo revolution. *Bulletin of the American Mathematical*  
148 *Society*, 46(2):179–205, 2009.
- 149 A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian data*  
150 *analysis*, volume 2. CRC press Boca Raton, FL, 2014.
- 151 E. Jang, S. Gu, and B. Poole. Categorical reparameterization with Gumbel-softmax. *arXiv preprint*  
152 *arXiv:1611.01144*, 2016.
- 153 S. Kato, H. Kaplan, T. Schrödel, S. Skora, T. Lindsay, E. Yemini, S. Lockery, and M. Zimmer. Global  
154 brain dynamics embed the motor command sequence of *Caenorhabditis elegans*. *Cell*, 163(3):656 –  
155 669, 2015. ISSN 0092-8674.
- 156 D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*,  
157 2013.
- 158 R. G. Krishnan, U. Shalit, and D. Sontag. Deep Kalman filters. *arXiv preprint arXiv:1511.05121*,  
159 2015.
- 160 H. Kushner and G. Yin. Stochastic approximation algorithms for parallel and distributed processing.  
161 *Stochastics: An International Journal of Probability and Stochastic Processes*, 22(3-4):219–250,  
162 1987.
- 163 S. Linderman, M. Johnson, A. Miller, R. Adams, D. Blei, and L. Paninski. Bayesian learning and  
164 inference in recurrent switching linear dynamical systems. In *Artificial Intelligence and Statistics*,  
165 pages 914–922, 2017a.
- 166 S. W. Linderman, G. E. Mena, H. Cooper, L. Paninski, and J. P. Cunningham. Reparameterizing the  
167 Birkhoff polytope for variational permutation inference. *arXiv preprint arXiv:1710.09508*, 2017b.
- 168 R. Lints, Z. F. Altun, H. Weng, T. Stephney, G. Stephney, M. Volaski, and D. H. Hall. WormAtlas  
169 Update. 2005.
- 170 C. J. Maddison, A. Mnih, and Y. W. Teh. The Concrete distribution: A continuous relaxation of  
171 discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 172 J. P. Nguyen, F. B. Shipley, A. N. Linder, G. S. Plummer, M. Liu, S. U. Setru, J. W. Shaevitz, and A. M.  
173 Leifer. Whole-brain calcium imaging with cellular resolution in freely behaving *Caenorhabditis*  
174 *elegans*. *Proceedings of the National Academy of Sciences*, 113(8):E1074–E1081, 2016.
- 175 L. Paninski, Y. Ahmadian, D. G. Ferreira, S. Koyama, K. R. Rad, M. Vidne, J. Vogelstein, and W. Wu.  
176 A new look at state-space models for neural data. *Journal of computational neuroscience*, 29(1-2):  
177 107–126, 2010.
- 178 G. Papandreou and A. L. Yuille. Perturb-and-map random fields: Using discrete optimization to learn  
179 and sample from energy models. In *Computer Vision (ICCV), 2011 IEEE International Conference*  
180 *on*, pages 193–200. IEEE, 2011.
- 181 D. Soudry, S. Keshri, P. Stinson, M.-h. Oh, G. Iyengar, and L. Paninski. Efficient" shotgun" inference  
182 of neural connectivity from highly sub-sampled activity data. *PLoS computational biology*, 11(10):  
183 e1004464, 2015.

- 184 L. R. Varshney, B. L. Chen, E. Paniagua, D. H. Hall, and D. B. Chklovskii. Structural properties of  
185 the *Caenorhabditis elegans* neuronal network. *PLoS Computational Biology*, 7(2):e1001066, 2011.
- 186 J. T. Vogelstein, J. M. Conroy, V. Lyzinski, L. J. Podrazik, S. G. Kratzer, E. T. Harley, D. E. Fishkind,  
187 R. J. Vogelstein, and C. E. Priebe. Fast approximate quadratic programming for graph matching.  
188 *PLOS one*, 10(4):e0121002, 2015.
- 189 J. G. White, E. Southgate, J. N. Thomson, and S. Brenner. The structure of the nervous system of  
190 the nematode *Caenorhabditis elegans*: the mind of a worm. *Phil. Trans. R. Soc. Lond*, 314:1–340,  
191 1986.
- 192 R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
193 learning. *Machine Learning*, 8(3–4):229–256, 1992.