

---

# Reparameterizing the Birkhoff Polytope for Variational Permutation Inference: Supplementary Material

---

Anonymous Authors  
Anonymous Institutions

## A Alternative methods of discrete variational inference

Recently there have been a number of proposals for extending the reparameterization trick [Rezende et al., 2014, Kingma and Welling, 2014] to high dimensional discrete problems<sup>1</sup> by relaxing them to analogous continuous problems [Maddison et al., 2016, Jang et al., 2016, Kusner and Hernández-Lobato, 2016]. These approaches are based on the following observation: if  $x \in \{0, 1\}^N$  is a one-hot vector drawn from a categorical distribution, then the support of  $p(x)$  is the set of vertices of the  $N - 1$  dimensional simplex. We can represent the distribution of  $x$  as an atomic density on the simplex.

### A.1 The Gumbel-softmax method

Viewing  $x$  as a vertex of the simplex motivates a natural relaxation: rather than restricting ourselves to atomic measures, consider continuous densities on the simplex. To be concrete, suppose the density of  $x$  is defined by the transformation,

$$\begin{aligned} \xi_n &\stackrel{\text{iid}}{\sim} \text{Gumbel}(0, 1) \\ \psi_n &= \log \theta_n + \xi_n \\ x &= \text{softmax}(\psi/\tau) \\ &= \left( \frac{e^{\psi_1/\tau}}{\sum_{n=1}^N e^{\psi_n/\tau}}, \dots, \frac{e^{\psi_N/\tau}}{\sum_{n=1}^N e^{\psi_n/\tau}} \right). \end{aligned}$$

The output  $x$  is now a point on the simplex, and the parameter  $\theta = (\theta_1, \dots, \theta_N)$  can be optimized via stochastic gradient ascent with the reparameterization trick.

The Gumbel distribution leads to a nicely interpretable model: when  $\theta$  is a probability mass function, adding Gumbel noise and taking the argmax yields an exact sample from  $\theta$ ; the softmax is a natural relaxation. As

---

<sup>1</sup>Discrete inference is only problematic in the high dimensional case, since in low dimensional problems we can enumerate the possible values of  $x$  and compute the normalizing constant  $p(y) = \sum_x p(y, x)$ .

the temperature  $\tau$  goes to zero, the softmax converges to the argmax function. Ultimately, however, this is just a continuous relaxation of an atomic density to a continuous density.

Stick-breaking and rounding offer two alternative ways of conceiving a relaxed version of a discrete random variable, and both are amenable to reparameterization. However, unlike the Gumbel-Softmax, these relaxations enable extensions to more complex combinatorial objects, notably, permutations.

### A.2 Stick-breaking

The stick-breaking transformation to the Birkhoff polytope presented in the main text contains a recipe for stick-breaking on the simplex. In particular, as we filled in the first row of the doubly-stochastic matrix, we were transforming a real-valued vector  $\psi \in \mathbb{R}^{N-1}$  to a point in the simplex. We present this procedure for discrete variational inference again here in simplified form. Start with a reparameterization of a Gaussian vector,

$$\begin{aligned} \xi_n &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \\ \psi_n &= \mu_n + \eta_n \xi_n, \quad 1 \leq n \leq N - 1, \end{aligned}$$

parameterized by  $\theta = (\mu_n, \eta_n)_{n=1}^{N-1}$ . Then map this to a point in the simplex:

$$\begin{aligned} x_1 &= \sigma(\psi_1/\tau), \\ x_n &= \sigma(\psi_n/\tau) \left( 1 - \sum_{m=1}^{n-1} x_m \right), \quad 2 \leq n \leq N - 1, \\ x_N &= 1 - \sum_{m=1}^{N-1} x_m, \end{aligned}$$

where  $\sigma(u) = (1 + e^{-u})^{-1}$  is the logistic function. Here,  $\sigma(\psi_n/\tau)$  is the fraction of the remaining “stick” of probability mass assigned to  $x_n$ . This procedure is invertible, the Jacobian  $\frac{dx}{d\psi}$  is lower-triangular, and the determinant of the Jacobian is easy to compute. [Linderman et al. \[2015\]](#) compute the density of  $x$  implied by a Gaussian density on  $\psi$ .

The temperature  $\tau$  controls how concentrated  $p(\pi)$  is at the vertices of the simplex, and with appropriate choices of parameters, in the limit  $\tau \rightarrow 0$  we can recover any categorical distribution. In the other limit, as  $\tau \rightarrow \infty$ , the density concentrates on a point in the interior of the simplex determined by the parameters, and for intermediate values, the density is continuous on the simplex.

Finally, note that the logistic-normal construction only one possible choice. We could instead let  $\psi_n \sim \text{Beta}(\frac{a_n}{\tau}, \frac{b_n}{\tau})$  and  $x_n = \psi_n$ . This would lead to the Dirichlet distribution on the simplex. The beta distribution is slightly harder to reparameterize since it is typically simulated with a rejection sampling procedure, but Naesseth et al. [2017] have shown how this can be handled with a mix of reparameterization and score-function gradients. Alternatively, the beta distribution could be replaced with the Kumaraswamy distribution, which is quite similar to the beta distribution but is easily reparameterizable.

### A.3 Rounding

Rounding transformations also have a natural analog for discrete variational inference. Define the rounding operator,

$$\text{round}(\psi) = \arg \min_{e_n} \|e_n - \psi\|^2,$$

which maps  $\psi \in \mathbb{R}^N$  to the one-hot vectors  $e_n$ ; i.e. the vectors in  $\{0, 1\}^N$  with  $n$ -th entry equal to one and all other entries equal zero. This is equivalent to defining  $\text{round}(\psi) = e_{n^*}$  where

$$\begin{aligned} n^* &= \arg \min_n \|e_n - \psi\|^2 \\ &= \arg \min_n \sum_{m \neq n} \psi_m^2 + (1 - \psi_n)^2 \\ &= \arg \min_n \sum_{m \neq n} \psi_m^2 + \psi_n^2 - 2\psi_n + 1 \\ &= \arg \min_n \|\psi\|^2 - 2\psi_n + 1 \\ &= \arg \max_n \psi_n. \end{aligned}$$

In the case of a tie, let  $n^*$  be the smallest index  $n$  such that  $\psi_n > \psi_m$  for all  $m < n$ . Rounding effectively partitions the space into  $N$  disjoint ‘‘Voronoi’’ cells,

$$V_n = \left\{ \psi \in \mathbb{R}^N : \psi_n \geq \psi_m \forall m \wedge \psi_n > \psi_m \forall m < n \right\}.$$

By definition,  $\text{round}(\psi) = e_{n^*}$  for all  $\psi \in V_{n^*}$ .

We define a map that pulls points toward their rounded values,

$$x = \tau\psi + (1 - \tau)\text{round}(\psi). \quad (1)$$

**Proposition 1.** *For  $\tau \in [0, 1]$ , the map defined by (1) moves points strictly closer to their rounded values so that  $\text{round}(\psi) = \text{round}(x)$ .*

*Proof.* Note that the Voronoi cells are intersections of halfspaces and, as such, are convex sets. Since  $x$  is a convex combination of  $\psi$  and  $e_{n^*}$ , both of which belong to the convex set  $V_{n^*}$ ,  $x$  must belong to  $V_{n^*}$  as well.  $\square$

Similarly,  $x$  will be a point on the simplex if and only if  $\psi$  is on the simplex as well. By analogy to the rounding transformations for permutation inference, in categorical inference we use a Gaussian distribution  $\psi \sim \mathcal{N}(\text{proj}(m), H)$ , where  $\text{proj}(m)$  is the projection of  $m \in \mathbb{R}_+^N$  onto the simplex. Still, the simplex has zero measure under the Gaussian distribution. It follows that the rounded points  $x$  will almost surely not be on the simplex either. The supposition of this approach is that this is not a problem: relaxing to the simplex is nice but not required.

In the zero-temperature limit we obtain a discrete distribution on the vertices of the simplex. For  $\tau \in (0, 1]$  we have a distribution on  $\mathcal{X}_\tau \subseteq \mathbb{R}^N$ , the subset of the reals to which the rounding operation maps. (For  $0 \leq \tau < 1$  this is a strict subset of  $\mathbb{R}^N$ .) To derive the density  $q(x)$ , we need the inverse transformation and the determinant of its Jacobian. From Proposition 1, it follows that the inverse transformation is given by,

$$\psi = \frac{1}{\tau}x - \frac{1 - \tau}{\tau}\text{round}(x).$$

As long as  $\psi$  is in the interior of its Voronoi cell, the  $\text{round}$  function is piecewise constant and the Jacobian is  $\frac{d\psi}{dx} = \frac{1}{\tau}I$ , and its determinant is  $\tau^{-N}$ . Taken together, we have,

$$\begin{aligned} q(x; m, H) &= \\ &\tau^{-N} \mathcal{N}\left(\frac{1}{\tau}x - \frac{1 - \tau}{\tau}\text{round}(x); \text{proj}(m), H\right) \\ &\quad \times \mathbb{I}[x \in \mathcal{X}_\tau]. \end{aligned}$$

Compare this to the density of the rounded random variables for permutation inference.

## B Limit analysis for stick-breaking

We show that stick-breaking for discrete variational inference can converge to any categorical distribution in the zero-temperature limit. We do so with a sequence of propositions: first we show that in the zero-temperature limit, the distribution of  $\sigma(\psi_n/\tau)$  converges to a Bernoulli distribution. Then we show that when  $\sigma(\psi_n/\tau)$  is Bernoulli (rather than a continuous

density on the unit interval), the distribution on  $x$  obtained by applying the stick-breaking transformation to  $\psi$  is categorical.

**Proposition 2.** *Let  $z = \sigma(\psi/\tau)$  with  $\psi \sim \mathcal{N}(\mu, \eta^2)$ . In the limit  $\tau \rightarrow 0$  we have  $z \sim \text{Bern}(\Phi(-\frac{\mu}{\eta}))$ , where  $\Phi(\cdot)$  denotes the Gaussian cumulative distribution function (cdf).*

*Proof.* Let  $F_z$  be the cdf of the random variable  $z$ . Since  $z$  is a random variable on the unit interval,  $F_z$  is a non-decreasing function on  $[0, 1]$  with  $F_z(0) = 0$  and  $F_z(1) = 1$ . Reparameterize  $\psi = \mu + \eta\xi$  where  $\xi \sim \mathcal{N}(0, 1)$ . Then we have,

$$\begin{aligned} F_z(u) &= \Pr(\sigma(\psi/\tau) < u) \\ &= \Pr(\psi < \tau\sigma^{-1}(u)) \\ &= \Pr(\xi < \frac{\tau}{\eta}\sigma^{-1}(u) - \frac{\mu}{\eta}) \\ &= \Phi(-\frac{\tau}{\eta}\sigma^{-1}(u) - \frac{\mu}{\eta}). \end{aligned}$$

By the continuity of  $\Phi$  we have,

$$\lim_{\tau \rightarrow 0} F_z(u) = \Phi(-\frac{\mu}{\eta}) \quad \text{for } u \in (0, 1).$$

This is the cdf of a Bernoulli random with probability  $\rho = \Phi(-\frac{\mu}{\eta})$ .  $\square$

**Proposition 3.** *As above, let  $z_n = \sigma(\psi_n/\tau)$ . When  $z_n \sim \text{Bern}(\rho_n)$  with  $\rho_n \in [0, 1]$  for  $n = 1, \dots, N$ , the random variable  $x$  obtained from applying the stick-breaking transformation to  $z$  will have an atomic distribution with atoms in the vertices of  $\Delta_N$ ; i.e.,  $x \sim \text{Cat}(\pi)$  where*

$$\begin{aligned} \pi_1 &= \rho_1 \\ \pi_n &= \rho_n \prod_{m=1}^{n-1} (1 - \rho_m) \quad n = 2, \dots, N-1, \\ \pi_N &= \prod_{m=1}^{N-1} (1 - \rho_m). \end{aligned}$$

*Proof.* From the stick-breaking definition,  $x_1 = z_1$ ,  $x_n = z_n(1 - \sum_{m < n} x_m)$ , and  $x_N = 1 - \sum_{m < N} x_m$ . When  $z_n \in \{0, 1\}$  for all  $n = 1, \dots, N-1$ , we have the following equivalencies. For the first element,

$$x_1 = 1 \iff z_1 = 1;$$

for  $1 < n < N-1$ :

$$x_n = 1 \iff (z_n = 1) \bigwedge_{m=1}^{n-1} (z_m = 0);$$

and for the last element,

$$x_N = 1 \iff \bigwedge_{m=1}^{N-1} (z_m = 0).$$

These events are mutually exclusive, implying that  $x$  will necessarily be a one-hot vector, i.e. a categorical random variable. Since  $z_1, \dots, z_{N-1}$  are independent Bernoulli random variables, the probabilities of these events are given by the  $\pi, \dots, \pi_N$  stated in the proposition.  $\square$

These two propositions, combined with the invertibility of the stick-breaking procedure, lead to our main result.

**Lemma 1.** *In the zero-temperature limit, stick-breaking of logistic-normal random variables can realize any categorical distribution on  $x$ .*

*Proof.* There is a one-to-one correspondence between  $\pi \in \Delta_N$  and  $\rho \in [0, 1]^{N-1}$ . Specifically,

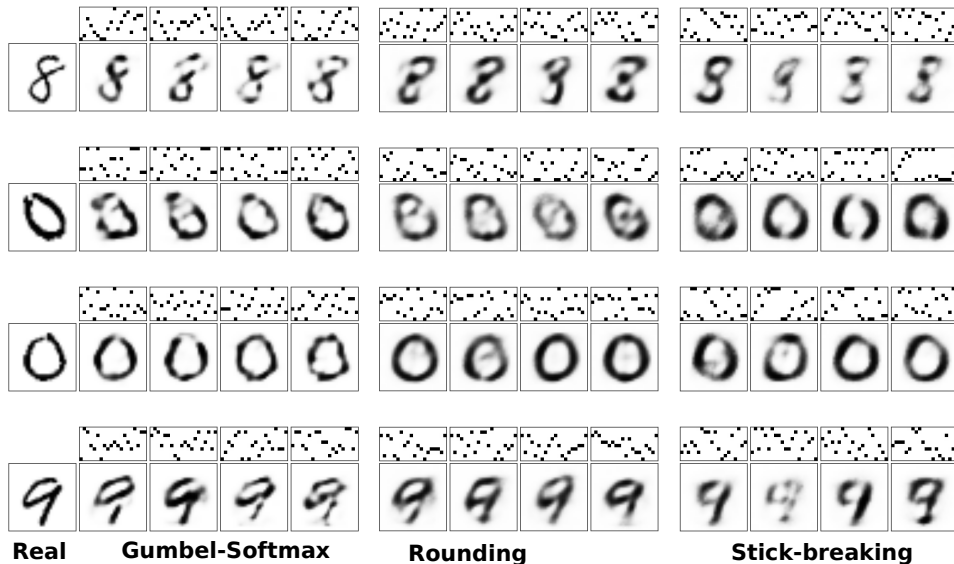
$$\begin{aligned} \rho_1 &= \pi_1 \\ \rho_n &= \frac{\pi_n}{\prod_{m=1}^{n-1} (1 - \rho_m)} \quad \text{for } n = 2, \dots, N-1. \end{aligned}$$

Since these are recursively defined, we can substitute the definition of  $\rho_m$  to obtain an expression for  $\rho_n$  in terms of  $\pi$  only. Thus, by Proposition 3, any desired categorical distribution  $\pi$  implies a set of Bernoulli parameters  $\rho$ . From Proposition 2, in the zero temperature limit, any desired  $\rho_n$  can be obtained with appropriate choice of Gaussian mean  $\mu_n$  and variance  $\eta_n^2$ . Thus, stick-breaking can realize any categorical distribution when  $\tau \rightarrow 0$ .  $\square$

## C Variational Autoencoders (VAE) with categorical latent variables

We considered the density estimation task on MNIST digits, as in Maddison et al. [2016], Jang et al. [2016], where observed digits are reconstructed from a latent discrete code. We used the continuous ELBO for training, and evaluated performance based on the marginal likelihood, estimated through the multi-sample variational objective of the discretized model. We compared against the methods of Jang et al. [2016], Maddison et al. [2016] and obtained the results in Table 1. While stick-breaking and rounding fare slightly worse than the Gumbel-softmax method, they are readily extensible to more complex discrete objects, as shown in the main paper.

Figure 1 shows MNIST reconstructions using Gumbel-Softmax, stick-breaking and rounding reparameterizations. In all the three cases reconstructions are reasonably accurate, and there is diversity in reconstructions.



**Figure 1:** Examples of true and reconstructed digits from their corresponding random codes using with  $K = 20$  categorical variables with  $N = 10$  possible values.

**Table 1:** Summary of results in VAE

Method	$-\log p(x)$
Gumbel-Softmax	106.7
Concrete	111.5
Rounding	121.1
Stick-breaking	119.8

sians for each dimension,

$$p(X) = \prod_{m=1}^N \prod_{n=1}^N \frac{1}{2} (\mathcal{N}(x_{mn} | 0, \eta^2) + \mathcal{N}(x_{mn} | 1, \eta^2)).$$

Although this prior puts significant mass around invalid points (e.g. **1**), it penalizes  $X$  that far from  $\mathcal{B}_N$ .

## D Variational permutation inference details

Here we discuss more of the subtleties of variational permutation inference and present the mathematical derivations in more detail.

### D.1 Continuous prior distributions.

Continuous relaxations require re-thinking the objective. As in [Maddison et al. \[2016\]](#), we maximize a relaxed ELBO, for which we need to specify a new continuous prior  $p(X)$  over the relaxed discrete latent variables, here, over relaxations of permutation matrices. Moreover, it is critical to design sensible priors for relaxed permutations. Ideally, this prior should penalize values of  $X$  that are far from permutation matrices.

For our categorical experiment on MNIST we use a mixture of Gaussians around each vertex,  $p(x) = \frac{1}{N} \sum_{n=1}^N \mathcal{N}(x | e_n, \eta^2)$ . This can be extended to permutations, where we use a mixture of Gaus-

### D.2 Deriving an expression for the ELBO

Here we show that if  $X = g(\Xi; \theta)$  with  $g$  differentiable one can evaluate the second term in equation (??). Moreover, both the stick-breaking and rounding transformations factor as  $g = h \circ f$  with  $X = h(\Psi)$  and  $\Psi = f(\Xi; \theta)$ . (Both  $h$  and  $f$  are invertible.) This means all dependency of  $X$  in the parameters is through the random variable  $\Psi$  with implicit density  $p(\Psi; \theta)$ .

We compute the entropy of  $q(X; \theta)$  (the second term in the ELBO) by,

$$\begin{aligned} \mathbb{E}_{q(X; \theta)} [-\log q(X; \theta)] &= \mathbb{E}_{r(\Xi)} [-\log q(h(f(\Xi, \theta)); \theta)] \\ &= \mathbb{E}_{p(\Psi; \theta)} [-\log q(h(\Psi); \theta)]. \end{aligned}$$

where the second equality follows by the “law of the unconscious statistician.”

Now, by the change of variable theorem and derivative

and determinant inversion rules,

$$\begin{aligned} q(X; \theta) &= p(h^{-1}(X); \theta) \left| \frac{dh^{-1}(X)}{dX} \right| \\ &= p(h^{-1}(X); \theta) \left| \frac{dh(\Psi)}{d\Psi} \right|_{\Psi=h^{-1}(X)}^{-1}. \end{aligned}$$

Now we appeal once more to the law of the unconscious statistician,

$$\begin{aligned} \mathbb{E}_{q(X; \theta)} [-\log q(X; \theta)] &= \mathbb{E}_{p(\Psi; \theta)} \left[ -\log p(\Psi; \theta) + \log \left| \frac{dh(\Psi)}{d\Psi} \right| \right] \\ &= \mathbb{H}(\Psi; \theta) + \mathbb{E}_{r(\Xi)} \left[ \left| \frac{dh(\Psi; \theta)}{d\Psi} \right|_{\Psi=f(\Xi; \theta)} \right]. \end{aligned}$$

### Estimating the ELBO

Here we describe how to compute each of terms of equation (??), needed for ELBO computations. First, as  $\Psi$  is Gaussian for both rounding and stick-breaking, the entropy term is straightforward and equal to  $N \log(\eta^2 2\pi e)/2$  ( $\eta$  may depend on the temperature and depends on the method).

Notice that to state  $\Psi$  is Gaussian in the stick-breaking case we slightly deviate from ???. Specifically, here we call  $\Psi = \frac{\mu_{mn} + \eta_{mn} \Xi_{mn}}{\tau}$  and define  $\Psi' = \sigma(\Psi)$ .

The second term of equation (??) is estimated using Monte-Carlo samples, and its derivation depends on the method.

### Rounding

Here  $H$  is piecewise linear: the set of discontinuities (border of the 'Voronoi cells' associated to each permutation) has Lebesgue measure zero. So we can still apply the change of variables theorem. Therefore,  $\log |DH(F(\Xi; \theta))| = N \log \tau$ . This means we don't even need to take samples to compute this term.

### Stick-breaking

It is important to note that the transformation  $H$  that maps  $\Psi' \rightarrow X$  is only piecewise continuous: the function is not differentiable at the points where the bounds change; for example, when changing  $\Psi'$  causes the active upper bound to switch from the row to the column constraint or vice versa.

Notice that these bounds only depend on values of  $X$  that have already been computed; i.e., those that are above or to the left of the  $(i, j)$ -th entry. Thus, the transformation from  $\Psi'$  to  $X$  is feed-forward according to this ordering. Consequently, the Jacobian of the inverse transformation  $H^{-1}$ ,  $d\Psi'/dX$ , is lower triangular, and its determinant is the product of its

diagonal,

$$\begin{aligned} \left| \frac{d\Psi'}{dX} \right| &= \prod_{i=1}^{n-1} \prod_{j=1}^{n-1} \frac{\partial \psi_{ij}}{\partial \pi_{ij}} \\ &= \prod_{i=1}^{n-1} \prod_{j=1}^{n-1} \frac{\partial}{\partial \pi_{ij}} \sigma^{-1} \left( \frac{\pi_{ij} - \ell_{ij}}{u_{ij} - \ell_{ij}} \right) \\ &= \prod_{i=1}^{n-1} \prod_{j=1}^{n-1} \left( \frac{1}{u_{ij} - \ell_{ij}} \right) \left( \frac{u_{ij} - \ell_{ij}}{\pi_{ij} - \ell_{ij}} \right) \left( \frac{u_{ij} - \ell_{ij}}{u_{ij} - \pi_{ij}} \right) \\ &= \prod_{i=1}^{n-1} \prod_{j=1}^{n-1} \frac{u_{ij} - \ell_{ij}}{(\pi_{ij} - \ell_{ij})(u_{ij} - \pi_{ij})} \end{aligned}$$

To compute the gradient of the forward transformation  $H$  one simply needs to invert the above (or put a negative sign, in the logarithm scale). Finally, to incorporate the effect of  $\sigma$  ( $\Psi' = \sigma(\Psi)$ ), by the chain rule, one only needs to add a term corresponding to this derivative,  $d\sigma(x)/dx = \sigma(x)\sigma(-x)$ .

### Experiment details

Experiments were run on a High Performance Computing (HPC) cluster, allowing the execution of hundreds of processes in parallel to efficiently determine best hyperparameter configurations.

For experiments with Variational Auto-encoder we used Tensorflow [Abadi et al., 2016], slightly changing the code made available in conjunction with ?. For experiments on synthetic matching and the C. elegans example we used Autograd [?], explicitly avoiding propagating gradients through the non-differentiable operation of solving a matching problem (the round in ??).

In all experiment we used the ADAM as optimizer, with learning rate 0.1. For rounding, the parameter vector  $H$  defined in ??(iii) was constrained to lie in the interval  $[0.1, 0.5]$ . Also, for rounding, we used ten iterations of the Sinkhorn-Knopp algorithm, to obtain points in the Birkhoff polytope. For stick-breaking the variances  $\nu$  defined in ?? was constrained between  $1e-8$  and 1.0. In either case, the temperature parameter was calibrated using a grid search.

In the C. elegans example we considered the symmetrized version of the adjacency matrix described in [Varshney et al., 2011] (i.e. we used  $A' = (A + A^\top)/2$ , and the matrix  $W$  was chosen antisymmetric, with entries sampled randomly with the sparsity pattern dictated by  $A'$ . To avoid divergence, the matrix  $W$  was then re-scaled by 1.1 times its spectral radius. This choice, although not essential, induced a reasonably well behaved linear dynamical system, rich in non-damped oscillations. We used a time window of

$T = 1000$  time samples, and added spherical standard noise at each time

## References

- M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- E. Jang, S. Gu, and B. Poole. Categorical reparameterization with Gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.
- M. J. Kusner and J. M. Hernández-Lobato. GANs for sequences of discrete elements with the Gumbel-softmax distribution. *arXiv preprint arXiv:1611.04051*, 2016.
- S. Linderman, M. Johnson, and R. P. Adams. Dependent multinomial models made easy: Stick-breaking with the polygamma augmentation. In *Advances in Neural Information Processing Systems*, pages 3456–3464, 2015.
- C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- C. Naesseth, F. Ruiz, S. Linderman, and D. Blei. Reparameterization gradients through acceptance-rejection sampling algorithms. In *Artificial Intelligence and Statistics*, pages 489–498, 2017.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of The 31st International Conference on Machine Learning*, pages 1278–1286, 2014.
- L. R. Varshney, B. L. Chen, E. Paniagua, D. H. Hall, and D. B. Chklovskii. Structural properties of the *Caenorhabditis elegans* neuronal network. *PLoS computational biology*, 7(2):e1001066, 2011.