

# Streaming, Distributed Variational Inference for Bayesian Nonparametrics

**Authored by:**

John W. Fisher III  
Trevor Campbell  
Jonathan P. How  
Julian Straub

## **Abstract**

This paper presents a methodology for creating streaming, distributed inference algorithms for Bayesian nonparametric (BNP) models. In the proposed framework, processing nodes receive a sequence of data minibatches, compute a variational posterior for each, and make asynchronous streaming updates to a central model. In contrast to previous algorithms, the proposed framework is truly streaming, distributed, asynchronous, learning-rate-free, and truncation-free. The key challenge in developing the framework, arising from fact that BNP models do not impose an inherent ordering on their components, is finding the correspondence between minibatch and central BNP posterior components before performing each update. To address this, the paper develops a combinatorial optimization problem over component correspondences, and provides an efficient solution technique. The paper concludes with an application of the methodology to the DP mixture model, with experimental results demonstrating its practical scalability and performance.

## **1 Paper Body**

Bayesian nonparametric (BNP) stochastic processes are streaming priors ? their unique feature is that they specify, in a probabilistic sense, that the complexity of a latent model should grow as the amount of observed data increases. This property captures common sense in many data analysis problems ? for example, one would expect to encounter far more topics in a document corpus after reading 106 documents than after reading 10 ? and becomes crucial in settings with unbounded, persistent streams of data. While their fixed, parametric cousins can be used to infer model complexity for datasets with known magnitude a priori [1, 2], such priors are silent with respect to notions of model complexity growth in streaming data settings. Bayesian nonparametrics are also

naturally suited to parallelization of data processing, due to the exchangeability, and thus conditional independence, they often exhibit via de Finetti's theorem. For example, labels from the Chinese Restaurant process [3] are rendered i.i.d. by conditioning on the underlying Dirichlet process (DP) random measure, and feature assignments from the Indian Buffet process [4] are rendered i.i.d. by conditioning on the underlying beta process (BP) random measure. Given these properties, one might expect there to be a wealth of inference algorithms for BNPs that address the challenges associated with parallelization and streaming. However, previous work has only addressed these two settings in concert for parametric models [5, 6], and only recently has each been addressed individually for BNPs. In the streaming setting, [7] and [8] developed streaming inference for DP mixture models using sequential variational approximation. Stochastic variational inference [9] and related methods [10?13] are often considered streaming algorithms, but their performance depends on the choice of a learning rate and on the dataset having known, fixed size a priori [5]. Outside of variational approaches, which are the focus of the present paper, there exist exact parallelized MCMC methods for BNPs [14, 15]; the tradeoff in using such methods is that they provide samples from the posterior rather than the distribution itself, and results regarding assessing 1

?m3 Minibatch Posterior  
?i2 Retrieve Original  
?o1  
Central Posterior  
Central Node  
Intermediate Posterior  
?o1  
?i1  
Node  
Central Node  
?3  
?  
?m3  
?3 (= ?m2 )  
?m2  
?i2  
?2  
?m2  
?2  
(= ?i2 + ?m3 ? ?0 )  
?m1  
?i1  
?1  
?m1  
?1  
(= ?i1 + ?m1 ? ?o1 )  
Node

```

Central Node
Node
Central Node
Node
Retrieve Data
...
yt+1
yt
...
Data Stream
(a) Retrieve the data/prior
...
yt+1
yt
...
...
Data Stream
yt+1
yt
...
...
yt+2
(c) Perform component ID
...
Data Stream
Data Stream
(b) Perform inference
yt+1
(d) Update the model

```

Figure 1: The four main steps of the algorithm that is run asynchronously on each processing node.

convergence remain limited. Sequential particle filters for inference have also been developed [16], but these suffer issues with particle degeneracy and exponential forgetting. The main challenge posed by the streaming, distributed setting for BNPs is the combinatorial problem of component identification. Most BNP models contain some notion of a countably infinite set of latent ?components? (e.g. clusters in a DP mixture model), and do not impose an inherent ordering on the components. Thus, in order to combine information about the components from multiple processors, the correspondence between components must first be found. Brute force search is in2 tractable even for moderately sized models ? there are  $K_1K+K$  possible correspondences for two 1 sets of components of sizes  $K_1$  and  $K_2$  . Furthermore, there does not yet exist a method to evaluate the quality of a component correspondence for BNP models. This issue has been studied before in the MCMC literature, where it is known as the ?label switching problem?, but past solution techniques are generally model-specific and restricted to use on very simple mixture models [17, 18]. This

paper presents a methodology for creating streaming, distributed inference algorithms for Bayesian nonparametric models. In the proposed framework (shown for a single node A in Figure 1), processing nodes receive a sequence of data minibatches, compute a variational posterior for each, and make asynchronous streaming updates to a central model using a mapping obtained from a component identification optimization. The key contributions of this work are as follows. First, we develop a minibatch posterior decomposition that motivates a learning-rate-free streaming, distributed framework suitable for Bayesian nonparametrics. Then, we derive the component identification optimization problem by maximizing the probability of a component matching. We show that the BNP prior regularizes model complexity in the optimization; an interesting side effect of this is that regardless of whether the minibatch variational inference scheme is truncated, the proposed algorithm is truncation-free. Finally, we provide an efficiently computable regularization bound for the Dirichlet process prior based on Jensen’s inequality<sup>1</sup>. The paper concludes with applications of the methodology to the DP mixture model, with experimental results demonstrating the scalability and performance of the method in practice.

2

## Streaming, distributed Bayesian nonparametric inference

The proposed framework, motivated by a posterior decomposition that will be discussed in Section 2.1, involves a collection of processing nodes with asynchronous access to a central variational posterior approximation (shown for a single node in Figure 1). Data is provided to each processing node as a sequence of minibatches. When a processing node receives a minibatch of data, it obtains the central posterior (Figure 1a), and using it as a prior, computes a minibatch variational posterior approximation (Figure 1b). When minibatch inference is complete, the node then performs component identification between the minibatch posterior and the current central posterior, accounting for possible modifications made by other processing nodes (Figure 1c). Finally, it merges the minibatch posterior into the central variational posterior (Figure 1d). In the following sections, we use the DP mixture [3] as a guiding example for the technical development of the inference framework. However, it is emphasized that the material in this paper generalizes to many other BNP models, such as the hierarchical DP (HDP) topic model [19], BP latent feature model [20], and Pitman-Yor (PY) mixture [21] (see the supplement for further details).<sup>1</sup>

Regularization bounds for other popular BNP priors may be found in the supplement.

2

### 2.1

#### Posterior decomposition

Consider a DP mixture model [3], with cluster parameters  $\theta$ , assignments  $z$ , and observed data  $y$ . For each asynchronous update made by each processing node, the dataset is split into three subsets  $y = y_0 \cup y_i \cup y_m$  for analysis. When the processing node receives a minibatch of data  $y_m$ , it queries the central processing node for the original posterior  $p(\theta, z_0 | y_0)$ , which will be used as the prior for minibatch inference. Once inference is complete, it

again queries the central processing node for the intermediate posterior  $p(\theta, z_0, z_i - y_0, y_i)$  which accounts for asynchronous updates from other processing nodes since minibatch inference began. Each subset  $y_r, r \in \{o, i, m\}$ , has  $N_r$  observations  $\{y_{rj}\}_{j=1}^N$ , and each variable  $z_{rj} \in \mathcal{N}$  assigns  $y_{rj}$  to cluster parameter  $\theta_{z_{rj}}$ . Given the independence of  $\theta$  and  $z$  in the prior, and the conditional independence of the data given the latent parameters, Bayes' rule yields the following decomposition of the posterior of  $\theta$  and  $z$  given  $y$ , Updated Central Posterior

$$\begin{aligned} z &\sim \{p(\theta, z - y) \theta\} \\ \text{Original Posterior} \\ \text{Minibatch Posterior} \\ \text{Intermediate Posterior} \\ z &\sim \{z\} \sim \{z\} \sim \{p(z_i, z_m - z_0) \theta p(\theta, z_0 - y_0) \theta p(\theta, z_m, z_0 - y_m, y_0) \theta p(\theta, z_i, z_0 - y_i, y_0) \cdot p(z_i - z_0) p(z_m - z_0)\} \\ (1) \end{aligned}$$

This decomposition suggests a simple streaming, distributed, asynchronous update rule for a processing node: first, obtain the current central posterior density  $p(\theta, z_0 - y_0)$ , and using it as a prior, compute the minibatch posterior  $p(\theta, z_m, z_0 - y_0, y_m)$ ; and then update the central posterior density by using (1) with the current central posterior density  $p(\theta, z_i, z_0 - y_i, y_0)$ . However, there are two issues preventing the direct application of the decomposition rule (1): Unknown component correspondence: Since it is generally intractable to find the minibatch posteriors  $p(\theta, z_m, z_0 - y_0, y_m)$  exactly, approximate methods are required. Further, as (1) requires the multiplication of densities, sampling-based methods are difficult to use, suggesting a variational approach. Typical mean-field variational techniques introduce an artificial ordering of the parameters in the posterior, thereby breaking symmetry that is crucial to combining posteriors correctly using density multiplication [6]. The use of (1) with mean-field variational approximations thus requires first solving a component identification problem. Unknown model size: While previous posterior merging procedures required a 1-to-1 matching between the components of the minibatch posterior and central posterior [5, 6], Bayesian nonparametric posteriors break this assumption. Indeed, the datasets  $y_o, y_i$ , and  $y_m$  from the same nonparametric mixture model can be generated by the same, disjoint, or an overlapping set of cluster parameters. In other words, the global number of unique posterior components cannot be determined until the component identification problem is solved and the minibatch posterior is merged.

## Variational component identification

Suppose we have the following mean-field exponential family prior and approximate variational posterior densities in the minibatch decomposition (1),

$$\begin{aligned} p(\theta_k) &= h(\theta_k) e^{\theta_k^T A(\theta_k)} \\ q(\theta_k, z_0) &= h(\theta_k, z_0) e^{\theta_k^T A(\theta_k, z_0)} \end{aligned}$$

$$\begin{aligned}
& \prod_{k=1}^K p(\mathbf{z}_k, \mathbf{z}_m, \mathbf{z}_o \mid \mathbf{y}_m, \mathbf{y}_o) \prod_{i=1}^{K_i} q_i(\mathbf{z}_i, \mathbf{z}_o) = \prod_{m=1}^{K_m} p(\mathbf{z}_m) \prod_{o=1}^{K_o} q_o(\mathbf{z}_o) \\
& \prod_{k=1}^K h(\mathbf{z}_k) e^{\mathbf{z}_k^T \mathbf{T}(\mathbf{z}_k)} A(\mathbf{z}_k) \\
& \prod_{k=1}^K p(\mathbf{z}_k, \mathbf{z}_i, \mathbf{z}_o \mid \mathbf{y}_i, \mathbf{y}_o) \prod_{i=1}^{K_i} q_i(\mathbf{z}_i, \mathbf{z}_o) = \prod_{i=1}^{K_i} p(\mathbf{z}_i) \prod_{o=1}^{K_o} q_o(\mathbf{z}_o) \\
& \prod_{k=1}^K h(\mathbf{z}_k) e^{\mathbf{z}_k^T \mathbf{T}(\mathbf{z}_k)} A(\mathbf{z}_k) ,
\end{aligned}
\tag{2}$$

where  $\mathbf{r}(\mathbf{z})$ ,  $\mathbf{r} \in \{o, i, m\}$  are products of categorical distributions for the cluster labels  $\mathbf{z}_r$ , and the goal is to use the posterior decomposition (1) to find the updated posterior approximation  $K \ Y$

$$\begin{aligned}
& p(\mathbf{z}, \mathbf{z} \mid \mathbf{y}) \prod_{k=1}^K q(\mathbf{z}, \mathbf{z}) = \prod_{k=1}^K p(\mathbf{z}_k) \\
& \prod_{k=1}^K h(\mathbf{z}_k) e^{\mathbf{z}_k^T \mathbf{T}(\mathbf{z}_k)} A(\mathbf{z}_k) .
\end{aligned}
\tag{3}$$

As mentioned in the previous section, the artificial ordering of components causes the naive application of (1) with variational approximations to fail, as disparate components from the approximate posteriors may be merged erroneously. This is demonstrated in Figure 3a, which shows results from a synthetic experiment (described in Section 4) ignoring component identification. As the number of parallel threads increases, more matching mistakes are made, leading to decreasing model quality. To address this, first note that there is no issue with the first  $K_o$  components of  $q_m$  and  $q_i$ ; these can be merged directly since they each correspond to the  $K_o$  components of  $q_o$ . Thus, the component 0 identification problem reduces to finding the correspondence between the last  $K_m = K_m - K_o$  components of the minibatch posterior and the last  $K_i = K_i - K_o$  components of the intermediate posterior. For notational simplicity (and without loss of generality), fix the component ordering 0 of the intermediate posterior  $q_i$ , and define  $\mathbf{z} : [K_m] \rightarrow [K_i + K_m]$  to be the 1-to-1 mapping from minibatch posterior component  $k$  to updated central posterior component  $\mathbf{z}(k)$ , where  $[K] := \{1, \dots, K\}$ . The fact that the first  $K_o$  components have no ordering ambiguity can be expressed as  $\mathbf{z}(k) = k - K_o \ \forall k \in [K_o]$ . Note that the maximum number of components after merging is  $K_i + K_m$ , 0 since each of the last  $K_m$  components in the minibatch posterior may correspond to new components in the intermediate posterior. After substituting the three variational approximations (2) into (1), the goal of the component identification optimization is to find the 1-to-1 mapping  $\mathbf{z}$  that yields the largest updated posterior normalizing constant, i.e. matches components with similar densities,  $\arg\max_{\mathbf{z}} \prod_{k=1}^{K_m} p(\mathbf{z}_k, \mathbf{z}_m \mid \mathbf{z}_o) \prod_{i=1}^{K_i} q_i(\mathbf{z}_i, \mathbf{z}_o) \prod_{o=1}^{K_o} q_o(\mathbf{z}_o) \prod_{k=1}^{K_m} h(\mathbf{z}_k) e^{\mathbf{z}_k^T \mathbf{T}(\mathbf{z}_k)} A(\mathbf{z}_k)$ .

$$\begin{aligned} & \sum_{k=1}^K q_m(z, z_m) = \sum_{k=1}^K q_m(z_m) \\ & K \text{ m } Y \\ & T \\ & h(\sum_{k=1}^K q_m(z, z_m) e^{\sum_{k=1}^K q_m(z, z_m) T} A(\sum_{k=1}^K q_m(z, z_m) T) \end{aligned} \quad (4)$$

$q(k) = k, q_k \in [K_0], q \text{ 1-to-1 } q \text{ } (z_{mj} = q(k)) = P(z_{mj} = k). (z_m)$  is the distribution such that  $P(z_m)$  Taking the where  $\sum_{m=1}^M$  logarithm of the objective and exploiting the mean-field decoupling allows the separation of the objective into a sum of two terms: one expressing the quality of the matching between components (the integral over  $z$ ), and one that regularizes the final model size (the sum over  $z$ ). While the first term is available in closed form, the second is in general not. Therefore, using the concavity of the logarithm function, Jensen's inequality yields a lower bound that can be used in place of the intractable original objective, resulting in the final component identification optimization:  $0 K_i + K_m$

$$\begin{aligned} & \sum_{k=1}^K \sum_{i=1}^I \arg \max_{z_i} \sum_{k=1}^K \sum_{i=1}^I \log p(z_i, z_m, z_o) \end{aligned} \quad (5)$$

s.t.  $\sum_{k=1}^K q_k = \sum_{i=1}^I q_i + \sum_{m=1}^M q_m \text{ } q(k) = k \text{ } q_k \in [K_0], q \text{ 1-to-1}.$

A more detailed derivation of the optimization may be found in the supplement.  $E$  denotes expectation under the distribution  $q_o(z_o)$ ,  $q_i(z_i)$ ,  $q_m(z_m)$ , and

$\sum_{k=1}^K q_k \in ([K_m]) \text{ } r_k \in K_r \text{ } \sum_{k=1}^K q_k = \sum_{r \in \{o, i, m\}} \sum_{k=1}^K q_k = \sum_{k=1}^K q_k \text{ } / \in ([K_m])$  where  $\in ([K_m])$  denotes the range of the mapping  $\in$ . The definitions in (6) ensure that the prior  $q_0$  is used whenever a posterior  $r \in \{i, m, o\}$  does not contain a particular component  $k$ . The intuition for the optimization (5) is that it combines finding component correspondences with high similarity (via the log-partition function) with a regularization term<sup>2</sup> on the final updated posterior model size. Despite its motivation from the Dirichlet process mixture, the component identification optimization (5) is not specific to this model. Indeed, the derivation did not rely on any properties specific to the Dirichlet process mixture; the optimization applies to any Bayesian nonparametric model with a set of  $\in$  components<sup>2</sup>, and a set of combinatorial  $\in$  indicators<sup>2</sup>  $z$ . For example, the optimization applies to the hierarchical Dirichlet process topic model [10] with topic word distributions  $\in$  and local-to-global topic correspondences  $z$ , and to the beta process latent feature model [4] with features  $\in$  and  $z$

h i

$\in$  This is equivalent to the KL-divergence regularization  $\in_{KL} q_o(z_o) \in q_i(z_i) \in q_m(z_m) p(z_i, z_m, z_o)$ .

4

binary assignment vectors  $z$ . The form of the objective in the component identification optimization (5) reflects this generality. In order to apply the proposed streaming, distributed method to a particular model, one simply needs a black-box variational inference algorithm that computes posteriors of the form (2), and a way to compute or bound the expectation in the objective of (5).

### 2.3 Updating the central posterior

To update the central posterior, the node first locks it and solves for  $q$  via (5). Locking prevents other nodes from solving (5) or modifying the central posterior, but does not prevent other nodes from reading the central posterior, obtaining minibatches, or performing inference; the synthetic experiment in Section 4 shows that this does not incur a significant time penalty in practice. Then the processing node transmits  $q$  and its minibatch variational posterior to the central processing node where the product decomposition (1) is used to find the updated central variational posterior  $q$  in (3), with parameters

$$K = \max_i K_i, \quad \mu(z) = \frac{1}{K} \sum_{i=1}^K \mu_i(z_i), \quad \sigma(z) = \frac{1}{K} \sum_{i=1}^K \sigma_i(z_i), \quad \text{for } z \in \mathcal{Z}.$$

Finally, the node unlocks the central posterior, and the next processing node to receive a new minibatch will use the above  $K$ ,  $\mu(z)$ , and  $\sigma(z)$  from the central node as their  $K_o$ ,  $\mu_o(z_o)$ , and  $\sigma_o(z_o)$ .

### 3

#### Application to the Dirichlet process mixture model

The expectation in the objective of (5) is typically intractable to compute in closed-form; therefore, a suitable lower bound may be used in its place. This section presents such a bound for the Dirichlet process, and discusses the application of the proposed inference framework to the Dirichlet process mixture model using the developed bound. Crucially, the lower bound decomposes such that the optimization (5) becomes a maximum-weight bipartite matching problem. Such problems are solvable in polynomial time [22] by the Hungarian algorithm, leading to a tractable component identification step in the proposed streaming, distributed framework.

#### 3.1 Regularization lower bound

For the Dirichlet process with concentration parameter  $\alpha > 0$ ,  $p(z_i, z_m, z_o)$  is the Exchangeable Partition Probability Function (EPPF) [23]  $Y p(z_i, z_m, z_o) = \frac{1}{K!} \sum_{\pi \in \Pi_K} \frac{K!}{n_\pi!} \prod_{k \in K} \frac{\alpha^{n_k}}{n_k!}$ , (8)

where  $n_k$  is the amount of data assigned to cluster  $k$ , and  $K$  is the set of labels of nonempty clusters. Given that the variational distribution  $q_r(z_r)$ ,  $r \in \{i, m, o\}$  is a product of independent categorical distributions  $q_r(z_r) = \prod_{j=1}^J q_{rj}(z_{rj})$ , Jensen's inequality may be used to bound the  $k=1$  regularization in (5) below (see the supplement for further details) by  $\sum_{i=1}^I K_i + \sum_{m=1}^M K_m$

$$E_{q_r} [\log p(z_i, z_m, z_o)] \geq$$

$\sum_{k=1}^K$

$$\frac{1}{n_k} \sum_{r \in \{i, m, o\}} \log q_{rk}(z_{rk}) + \log \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{r \in \{i, m, o\}} \log q_{rk}(z_{rk}) + C$$

$k=1$

$$s_{rk} = s_{ik} + s_{mk} + s_{ok},$$

(9)

$$t_{rk} = t_{ik} + t_{mk} + t_{ok},$$



where  $C$  is a constant with respect to the component mapping  $\mathbf{r}$ , and  $\mathbf{P} \mathbf{N} \mathbf{r} \mathbf{r} \mathbf{k} \mathbf{K} \mathbf{r} \mathbf{k} \mathbf{K} \mathbf{r} \mathbf{j} = 1 \log(1 \mathbf{r} \mathbf{r} \mathbf{j} \mathbf{k}) \mathbf{j} = 1 \mathbf{r} \mathbf{j} \mathbf{k} \mathbf{s} \mathbf{r} \mathbf{k} = \mathbf{r} \mathbf{r} \mathbf{r} \mathbf{o}, \mathbf{i}, \mathbf{m} \mathbf{t} \mathbf{r} \mathbf{k} = \mathbf{r} \mathbf{r} \mathbf{r} \mathbf{o}, \mathbf{i}, \mathbf{m}$   
 $0 \mathbf{k} \mathbf{l} \mathbf{K} \mathbf{r} 0 \mathbf{k} \mathbf{l} \mathbf{K} \mathbf{r} \mathbf{P} \mathbf{N} \mathbf{P} \mathbf{N} (10) \mathbf{m} \mathbf{m} \mathbf{k} \mathbf{r} \mathbf{r} ([\mathbf{K} \mathbf{m} \mathbf{j}]) \mathbf{k} \mathbf{r} \mathbf{r} ([\mathbf{K} \mathbf{m} \mathbf{j}]) \mathbf{j} = 1 \log(1 \mathbf{r} \mathbf{r} \mathbf{m} \mathbf{j} \mathbf{r} \mathbf{r} 1 (\mathbf{k})$   
 $) \mathbf{j} = 1 \mathbf{r} \mathbf{m} \mathbf{j} \mathbf{r} \mathbf{r} 1 (\mathbf{k}) \mathbf{s} \mathbf{r} \mathbf{r} \mathbf{m} \mathbf{k} = \mathbf{t} \mathbf{r} \mathbf{r} \mathbf{m} \mathbf{k} = . 0$

$$\mathbf{k} \mathbf{r} \mathbf{r} ([\mathbf{K} / \mathbf{m} \mathbf{j}])$$

$$\mathbf{k} \mathbf{r} \mathbf{r} ([\mathbf{K} / \mathbf{m} \mathbf{j}])$$

$$0$$

Note that the bound (9) allows incremental updates: after finding the optimal mapping  $\mathbf{r} \mathbf{r} \mathbf{r}$ , the central update (7) can be augmented by updating the values of  $\mathbf{s} \mathbf{k}$  and  $\mathbf{t} \mathbf{k}$  on the central node to  $\mathbf{r}$

$$\mathbf{r}$$

$$\mathbf{s} \mathbf{k} \mathbf{r} \mathbf{s} \mathbf{r} \mathbf{i} \mathbf{k} + \mathbf{s} \mathbf{r} \mathbf{r} \mathbf{m} \mathbf{k} + \mathbf{s} \mathbf{r} \mathbf{o} \mathbf{k} ,$$

$$\mathbf{t} \mathbf{k} \mathbf{r} \mathbf{t} \mathbf{r} \mathbf{i} \mathbf{k} + \mathbf{t} \mathbf{r} \mathbf{r} \mathbf{m} \mathbf{k} + \mathbf{t} \mathbf{r} \mathbf{o} \mathbf{k} . 5$$

$$(11)$$

$$600$$

$$140 \text{ Truth Lower Bound}$$

$$\text{Truth Lower Bound } 120$$

$$200$$

$$\text{Regularization}$$

$$\text{Regularization}$$

$$400$$

$$\text{Increasing } \mathbf{r}$$

$$0$$

$$-200$$

$$\text{Increasing } \mathbf{r} \text{ } 80$$

$$60$$

$$-400$$

$$-600$$

$$100$$

$$0$$

$$20$$

$$40$$

$$60$$

$$80$$

$$40 \text{ } 0.001 \text{ Certain}$$

$$100$$

$$\text{Number of Clusters}$$

$$0.01$$

$$(a)$$

$$0.1$$

$$\text{Clustering Uncertainty } \mathbf{r}$$

$$1.0$$

$$1000 \text{ Uncertain}$$

$$(b)$$

Figure 2: The Dirichlet process regularization and lower bound, with (2a) fully uncertain labelling and varying number of clusters, and (2b) the number of clusters fixed with varying labelling uncertainty.

As with  $K$ ,  $\tau_k$ , and  $\tau$  from (7), after performing the regularization statistics update (11), a processing node that receives a new minibatch will use the above  $s_k$  and  $t_k$  as their  $s_{ok}$  and  $t_{ok}$ , respectively. Figure 2 demonstrates the behavior of the lower bound in a synthetic experiment with  $N = 100$

datapoints for various DP concentration parameter values  $\alpha = 10^{-3}, 10^{-2}$ . The true regularization  $\log E[p(z)]$  was computed by sample approximation with 104 samples. In Figure 2a, the number of clusters  $K$  was varied, with symmetric categorical label weights set to  $K$  two important phenomena. First, the bound increases as  $K \rightarrow 0$ ; in other words, it gives preference to fewer, larger clusters, which is the typical BNP ‘rich get richer’ property. Second, the behavior of the bound as  $K \rightarrow N$  depends on the concentration parameter  $\alpha$  as  $\alpha$  increases, more clusters are preferred. In Figure 2b, the number of clusters  $K$  was fixed to 10, and the categorical

label weights were sampled from a symmetric Dirichlet distribution with parameter  $\alpha = 10^{-3}, 10^{-2}$ . This figure demonstrates that the bound does not degrade significantly with high labelling uncertainty, and is nearly exact for low labelling uncertainty. Overall, Figure 2a demonstrates that the proposed lower bound exhibits similar behaviors to the true regularization, supporting its use in the optimization (5).

### 3.2 Solving the component identification optimization

Given that both the regularization (9) and component matching score in the objective (5) decompose 0-1, the objective can be rewritten using a matrix of matching as a sum of terms for each  $k \in [K_i + K_m]$   $K_i + K_m$  scores  $R = R$  and selector variables  $X \in \{0, 1\}^{(K_i + K_m) \times (K_i + K_m)}$ . Setting  $X_{kj} = 1$  indicates that component  $k$  in the minibatch posterior is matched to component  $j$  in the intermediate posterior (i.e.  $\tau(k) = j$ ), providing a score  $R_{kj}$  defined using (6) and (10) as

$R_{kj} = A(\tau_{ij} + \tau_{mk} + \tau_{oj}) + 1 - \exp(-\tau_{ij} - \tau_{mk} - \tau_{oj}) \log \tau_{ij} + \log \tau_{mk} + \log \tau_{oj}$  (12) The optimization (5) can be rewritten in terms of  $X$  and  $R$  as

$$\begin{aligned} X & \in \arg\max \text{tr } X^T R X \\ X^T \mathbf{1} &= \mathbf{1}, X_{kk} = 1, \tau_k \in [K_o] \quad 0 \leq T X \in \{0, 1\}^{(K_i + K_m) \times (K_i + K_m)}, \mathbf{1} \\ &= [1, \dots, 1] \\ \text{s.t. } X \mathbf{1} &= \mathbf{1}, \end{aligned} \quad (13)$$

The first two constraints express the 1-to-1 property of  $\tau(\cdot)$ . The constraint  $X_{kk} = 1$  for  $k \in [K_o]$  fixes the upper  $K_o \times K_o$  block of  $X$  to  $I$  (due to the fact that the first  $K_o$  components are matched directly),  $0 \leq T X$  and the off-diagonal blocks to 0. Denoting  $X_0, R_0$  to be the lower right  $(K_i + K_m) \times (K_i + K_m)$  blocks of  $X, R$ , the remaining optimization problem is a linear assignment problem on  $X$  with cost matrix  $\tau R_0$ , which can be solved using the Hungarian algorithm [3]. Note that if  $K_m = K_o$  or  $K_i = K_o$ , this implies that no matching problem needs to be solved for the first  $K_o$  components of the minibatch posterior are

matched directly, and the last  $K_m$  are set as new components. In practical implementation of the framework, new clusters are typically discovered at a diminishing rate as more data are observed, so the number of matching problems that are solved likewise tapers off. The final optimal component mapping  $\pi^*$  is found by finding the nonzero elements of  $X^* : \pi^* = \arg\max_{\pi} \sum_{k \in [K_m]} X^*_{kj} \pi_k$

For the experiments in this work, we used the implementation at [github.com/hrldcpr/hungarian](https://github.com/hrldcpr/hungarian).

CPU Time (s)

- 5.0
- 12.0
- 6.0
- 10.0
- 7.0
- 6.0
- 8.0
- 4.0
- 9.0
- 10.0
- 2.0
- 10.0
- 11.0
- 0.0
- 7.0
- 6.0
- 8.0
- 4.0
- 9.0
- 2.0
- 0.0
- 1
- 2
- 4
- 8
- 16
- 24
- 32
- 40
- 48
- 1
- 2
- 4
- 8
- 16

24  
 32  
 40  
 48  
 SDA-DP Batch SVA SVI moVB SC  
 -7  
 -8  
 -9  
 -10  
 -11  
 -11.0  
 -5  
 -4  
 -3  
 -2  
 -1  
 (a) SDA-DP without component ID  
 (b) SDA-DP with component ID  
 45  
 # Clusters # Matchings True # Clusters  
 120  
 20  
 3  
 4  
 100  
 100  
 80  
 80  
 60  
 # Clusters # Matchings True # Clusters  
 120  
 Count  
 Count  
 25  
 2  
 140  
 35 30  
 1  
 (c) Test log-likelihood traces  
 140  
 40  
 0  
 Time (Log10 s)  
 # Threads  
 # Threads  
 Count

-6.0  
 8.0  
 8.0  
 -6  
 -5.0 CPU Time Test Log Likelihood  
 Test Log Likelihood Test Log Likelihood  
 CPU Time Test Log Likelihood  
 10.0  
 Test Log Likelihood CPU Time (s)  
 12.0  
 60  
 15  
 40  
 40  
 10  
 20  
 5 0  
 0  
 10  
 20  
 30  
 40  
 50  
 60  
 Merge Time (microseconds)  
 (d) CPU time for component ID  
 70  
 0  
 20  
 0  
 20  
 40  
 60  
 80  
 # Minibatches Merged  
 (e) Cluster/component ID counts  
 100  
 0  
 1  
 2  
 4  
 8  
 16  
 24  
 32  
 40

Figure 3: Synthetic results over 30 trials. (3a-3b) Computation time and test log likelihood for SDA-DP with varying numbers of parallel threads, with component identification disabled (3a) and enabled (3b). (3c) Test log likelihood traces for SDA-DP (40 threads) and the comparison algorithms. (3d) Histogram of computation time (in microseconds) to solve the component identification optimization. (3e) Number of clusters and number of component identification problems solved as a function of the number of minibatch updates (40 threads). (3f) Final number of clusters and matchings solved with varying numbers of parallel threads.

In this section, the proposed inference framework is evaluated on the DP Gaussian mixture with a normal-inverse-Wishart (NIW) prior. We compare the streaming, distributed procedure coupled with standard variational inference [24] (SDA-DP) to five state-of-the-art inference algorithms: memoized online variational inference (moVB) [13], stochastic online variational inference (SVI) [9] with 1 learning rate  $(t+10)^{-2}$ , sequential variational approximation (SVA) [7] with cluster creation threshold 10 and prune/merge threshold  $10^{-3}$ , sub-cluster splits MCMC (SC) [14], and batch variational inference (Batch) [24]. Priors were set by hand and all methods were initialized randomly. Methods that use multiple passes through the data (e.g. moVB, SVI) were allowed to do so. moVB was allowed to make birth/death moves, while SVI/Batch had fixed truncations. All experiments were performed on a computer with 24 CPU cores and 12GiB of RAM. Synthetic: This dataset consisted of 100,000 2-dimensional vectors generated from a Gaussian mixture model with 100 clusters and a NIW( $\mu_0, \Sigma_0, \nu_0, \Psi_0$ ) prior with  $\mu_0 = 0$ ,  $\Sigma_0 = 10^{-3}I$ ,  $\nu_0 = 1$ , and  $\Psi_0 = 4$ . The algorithms were given the true NIW prior, DP concentration  $\beta = 5$ , and minibatches of size 50. SDA-DP minibatch inference was truncated to  $K = 50$  components, and all other algorithms were truncated to  $K = 200$  components. Figure 3 shows the results from the experiment over 30 trials, which illustrate a number of important properties of SDA-DP. First and foremost, ignoring the component identification problem leads to decreasing model quality with increasing number of parallel threads, since more matching mistakes are made (Figure 3a). Second, if component identification is properly accounted for using the proposed optimization, increasing the number of parallel threads reduces execution time, but does not affect the final model quality (Figure 3b). Third, SDA-DP (with 40 threads) converges to the same final test log likelihood as the comparison algorithms in significantly reduced time (Figure 3c). Fourth, each component identification optimization typically takes  $\sim 10^{-5}$  seconds, and thus matching accounts for less than a millisecond of total computation and does not affect the overall computation time significantly (Figure 3d). Fifth, the majority of the component matching problems are solved within the first 80 minibatch updates (out of a total of 2,000) — afterwards, the true clusters have

all been discovered and the processing nodes contribute to those clusters rather than creating new ones, as per the discussion at the end of Section 3.2 (Figure 3e). Finally, increased parallelization can be advantageous in discovering the correct number of clusters; with only one thread, mistakes made early on are built upon and persist, whereas with more threads there are more component identification problems solved, and thus more chances to discover the correct clusters (Figure 3f). 7

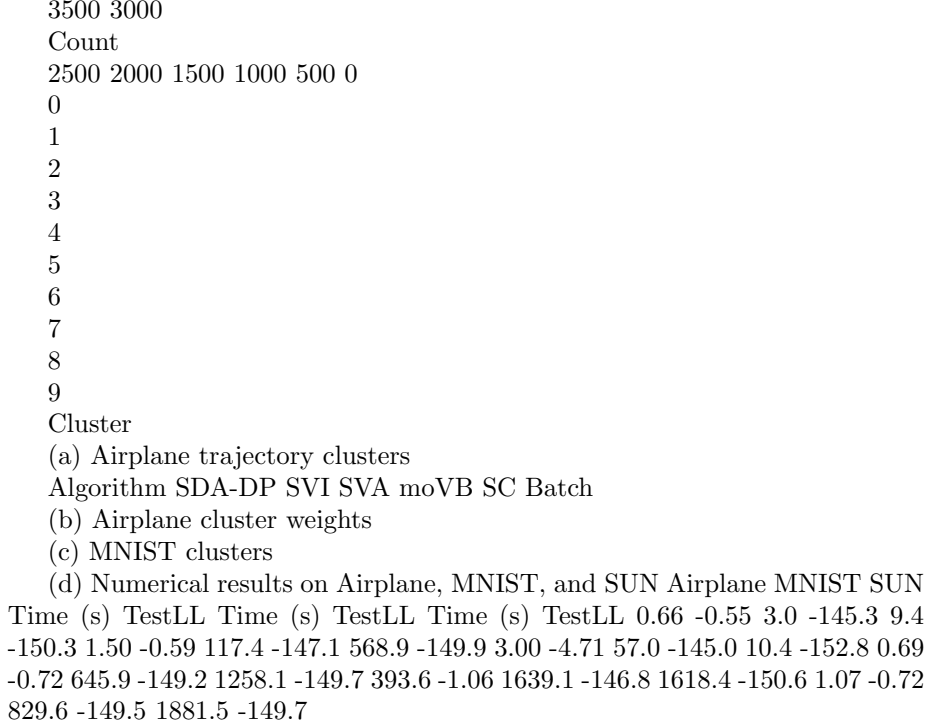


Figure 4: (4a-4b) Highest-probability instances and counts for 10 trajectory clusters generated by SDA-DP. (4c) Highest-probability instances for 20 clusters discovered by SDA-DP on MNIST. (4d) Numerical results.

**Airplane Trajectories:** This dataset consisted of 73,000,000 automatic dependent surveillance broadcast (ADS-B) messages collected from planes across the United States during the period 2013-03-22 01:30:00UTC to 2013-03-28 12:00:00UTC. The messages were connected based on plane call sign and time stamp, and erroneous trajectories were filtered based on reasonable spatial/temporal bounds, yielding 15,022 trajectories with 1,000 held out for testing. The latitude/longitude points in each trajectory were fit via linear regression, and the 3-dimensional parameter vectors were clustered. Data was split into minibatches of size 100, and SDA-DP used 16 parallel threads. **MNIST Digits [25]:** This dataset consisted of 70,000 28 x 28 images of hand-written digits, with 10,000 held out for testing. The images were reduced to 20 dimensions with PCA prior to clustering. Data was split into minibatches of size 500, and SDA-DP used 48 parallel threads. **SUN Images [26]:** This dataset consisted of 108,755 images

from 397 scene categories, with 8,755 held out for testing. The images were reduced to 20 dimensions with PCA prior to clustering. Data was split into minibatches of size 500, and SDA-DP used 48 parallel threads. Figure 4 shows the results from the experiments on the three real datasets. From a qualitative standpoint, SDA-DP discovers sensible clusters in the data, as demonstrated in Figures 4a-4c. However, an important quantitative result is highlighted by Table 4d: the larger a dataset is, the more the benefits of parallelism provided by SDA-DP become apparent. SDA-DP consistently provides a model quality that is competitive with the other algorithms, but requires orders of magnitude less computation time, corroborating similar findings on the synthetic dataset.

5

## Conclusions

This paper presented a streaming, distributed, asynchronous inference algorithm for Bayesian nonparametric models, with a focus on the combinatorial problem of matching minibatch posterior components to central posterior components during asynchronous updates. The main contributions are a component identification optimization based on a minibatch posterior decomposition, a tractable bound on the objective for the Dirichlet process mixture, and experiments demonstrating the performance of the methodology on large-scale datasets. While the present work focused on the DP mixture as a guiding example, it is not limited to this model – exploring the application of the proposed methodology to other BNP models is a potential area for future research. Acknowledgments This work was supported by the Office of Naval Research under ONR MURI grant N000141110688.

8

## 2 References

- [1] Agostino Nobile. Bayesian Analysis of Finite Mixture Distributions. PhD thesis, Carnegie Mellon University, 1994.
- [2] Jeffrey W. Miller and Matthew T. Harrison. A simple example of Dirichlet process mixture inconsistency for the number of components. In *Advances in Neural Information Processing Systems* 26, 2013.
- [3] Yee Whye Teh. Dirichlet processes. In *Encyclopedia of Machine Learning*. Springer, New York, 2010.
- [4] Thomas L. Griffiths and Zoubin Ghahramani. Infinite latent feature models and the Indian buffet process. In *Advances in Neural Information Processing Systems* 22, 2005.
- [5] Tamara Broderick, Nicholas Boyd, Andre Wibisono, Ashia C. Wilson, and Michael I. Jordan. Streaming variational Bayes. In *Advances in Neural Information Processing Systems* 26, 2013.
- [6] Trevor Campbell and Jonathan P. How. Approximate decentralized Bayesian inference. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*, 2014.
- [7] Dahua Lin. Online learning of nonparametric mixture models via sequential variational approximation. In *Advances in Neural Information Processing Systems* 26, 2013.
- [8] Xiaole Zhang, David J. Nott, Christopher Yau, and Ajay Jasra. A sequential algorithm for fast fitting of Dirichlet process mixture models. *Journal of Computational and Graphical*



Statistics, 23(4):1143–1162, 2014. [9] Matt Hoffman, David Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14:1303–1347, 2013. [10] Chong Wang, John Paisley, and David M. Blei. Online variational inference for the hierarchical Dirichlet process. In *Proceedings of the 11th International Conference on Artificial Intelligence and Statistics*, 2011. [11] Michael Bryant and Erik Sudderth. Truly nonparametric online variational inference for hierarchical Dirichlet processes. In *Advances in Neural Information Processing Systems* 23, 2009. [12] Chong Wang and David Blei. Truncation-free stochastic variational inference for Bayesian nonparametric models. In *Advances in Neural Information Processing Systems* 25, 2012. [13] Michael Hughes and Erik Sudderth. Memoized online variational inference for Dirichlet process mixture models. In *Advances in Neural Information Processing Systems* 26, 2013. [14] Jason Chang and John Fisher III. Parallel sampling of DP mixture models using sub-clusters splits. In *Advances in Neural Information Processing Systems* 26, 2013. [15] Willie Neiswanger, Chong Wang, and Eric P. Xing. Asymptotically exact, embarrassingly parallel MCMC. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*, 2014. [16] Carlos M. Carvalho, Hedibert F. Lopes, Nicholas G. Polson, and Matt A. Taddy. Particle learning for general mixtures. *Bayesian Analysis*, 5(4):709–740, 2010. [17] Matthew Stephens. Dealing with label switching in mixture models. *Journal of the Royal Statistical Society: Series B*, 62(4):795–809, 2000. [18] Ajay Jasra, Chris Holmes, and David Stephens. Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling. *Statistical Science*, 20(1):50–67, 2005. [19] Yee Whye Teh, Michael I. Jordan, Matthew J. Beal, and David M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006. [20] Finale Doshi-Velez and Zoubin Ghahramani. Accelerated sampling for the Indian buffet process. In *Proceedings of the International Conference on Machine Learning*, 2009. [21] Avinava Dubey, Sinead Williamson, and Eric Xing. Parallel Markov chain Monte Carlo for Pitman-Yor mixture models. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*, 2014. [22] Jack Edmonds and Richard Karp. Theoretical improvements in algorithmic efficiency for network flow problems. *Journal of the Association for Computing Machinery*, 19:248–264, 1972. [23] Jim Pitman. Exchangeable and partially exchangeable random partitions. *Probability Theory and Related Fields*, 102(2):145–158, 1995. [24] David M. Blei and Michael I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1(1):121–144, 2006. [25] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. MNIST database of handwritten digits. Online: [yann.lecun.com/exdb/mnist](http://yann.lecun.com/exdb/mnist). [26] Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. SUN 397 image database. Online: [vision.cs.princeton.edu/projects/2010/SUN](http://vision.cs.princeton.edu/projects/2010/SUN).