# An Infinite Hidden Markov Model With Local Transitions

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

We describe a generalization of the Hierarchical Dirichlet Process Hidden Markov Model (HDP-HMM) which is able to encode prior information that state transitions are more likely between "similar" states. This is accomplished by defining a similarity kernel on the state space, and scaling transition probabilities by pairwise similarities. This induces a global correlation structure over the transition probabilities based on the topology induced by the similarity kernel. We call this model the Hierarchical Dirichlet Process Hidden Markov Model with Local Transitions (HDP-HMM-LT). Unfortunately the conditional posterior of the transition distributions are no longer conjugate, due to the varying scale parameters, and so we present an alternative representation of this process as the marginalization of a Markov Jump Process in which: (1) some jump attempts fail, and (2) the probability of success is proportional to the similarity between the source and destination states. When holding times and failed transitions are reintroduced as auxiliary data, conditional conjugacy is restored, admitting exact Gibbs sampling. Even without the LT modification, conditioning on the holding times simplifies inference for the concentration parameters of the HDP, and allows immediate generalization to Semi-Markov dynamics without additional data augmentation. We evaluate the model and inference method on a collection of speaker diarization data sets in which speakers form conversational groups, so that there is prior dependence among chains, but most transitions involve only one or two chains at a time (transitions are local). Our model compares favorably to both the the HDP-H(S)MM and Binary Factorial HMM without suffering in performance when the data is generated directly from the comparison models.

## 1 Background

The conventional Hierarchical Dirichlet Process Hidden Markov Model (HDP-HMM) [5] is a prior distribution on the transition matrix of a Hidden Markov Model with a countably infinite state space. The rows of the infinite matrix are coupled through their dependence on a common, discrete base measure, itself drawn from a Dirichlet Process (DP). The hierarchical structure ensures that, despite the infinite state space, a common set of destination states will be reachable with high probability from each source state. The generative process for the HDP-HMM is the following:

Each of a countably infinite set of states, indexed by $j$, has parameters, $\theta_j$, drawn from a base measure, $H$. A top-level set of state weights, $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots)$, is drawn from a stick-breaking process (GEM) with concentration parameter $\gamma > 0$.

$$\theta_j \overset{i.i.d.}{\sim} H \qquad \boldsymbol{\beta} \sim \mathsf{GEM}(\gamma) \tag{1}$$

The actual transition distribution, $\boldsymbol{\pi}_j$, from state $j$, is drawn from a DP with concentration $\alpha$ and base measure $\boldsymbol{\beta}$:

$$\boldsymbol{\pi}_j \overset{i.i.d}{\sim} DP(\alpha\boldsymbol{\beta}) \qquad j = 1, 2, \ldots \tag{2}$$

The hidden state sequence is then generated according to the $\boldsymbol{\pi}_j$. Let $z_t$ be the index of the chain's state at time $t$. Then we have

$$z_t \mid z_{t-1}, \boldsymbol{\pi}_{z_{t-1}} \sim \boldsymbol{\pi}_{z_{t-1}} \qquad t = 1, 2, \ldots, T \tag{3}$$

where $T$ is the length of the data sequence. Finally, the emission distribution for state $j$ is a function of $\theta_j$, so that we have

$$y_t \mid z_t, \theta_{z_t} \sim F(\theta_{z_t}) \tag{4}$$

A shortcoming of this model is that the generative process does not take into account the fact that the set of source states is the same as the set of destination states: that is, the distribution $\boldsymbol{\pi}_j$ has an element which corresponds to state $j$. Put another way, there is no special treatment of the diagonal of the transition matrix, so that self-transitions are no more likely *a priori* than transitions to any other state. The Sticky HDP-HMM [1] addresses this issue by adding an extra mass at location $j$ to the base measure of the DP that generates $\boldsymbol{\pi}_j$. That is, (2) is replaced by

$$\boldsymbol{\pi}_j \sim DP(\alpha\boldsymbol{\beta} + \kappa\delta_j). \tag{5}$$

An alternative approach that treats self-transitions as special is the HDP Hidden Semi-Markov Model (HDP-HSMM) [3], wherein state duration distributions are modeled separately, and ordinary self-transitions are ruled out. In both the Sticky HDP-HMM and the HDP-HSMM, auxiliary latent variables are introduced to simplify conditional posterior distributions and facilitate Gibbs sampling. However, while both of these models have the ability to privilege self-transitions, they contain no notion of similarity for pairs of states that are not identical: in both cases, when the transition matrix is integrated out, the prior probability of transitioning to state $j'$ depends only on the top-level stick weight associated with state $j'$, and not on the identity or parameters of the previous state $j$.

## 2  An HDP-HMM With Local Transitions

We wish to add to the transition model the concept of a transition to a "nearby" state, where nearness of $j$ and $j'$ is a function of $\theta_j$ and $\theta_{j'}$. In order to accomplish this, we first consider an alternative construction of the transition distributions, based on the Normalized Gamma Process.

### 2.1  A Normalized Gamma Process representation of the HDP-HMM

We can define a random measure, $\mu = \sum_{j=1}^{\infty} \pi_j \delta_{\theta_j}$, where

$$\pi_j \overset{ind}{\sim} \mathsf{Gamma}(w_j, 1) \qquad T = \sum_{j=1}^{\infty} \pi_j \qquad \tilde{\pi}_j = \frac{\pi_j}{T} \tag{6}$$

and subject to the constraint that $\sum_{j \geq 1} w_j < \infty$. It follows [4] that $\mu$ is distributed as a Dirichlet Process with base measure $\mathbf{w} = \sum_{j=1}^{\infty} w_j \delta_{\theta_j}$. If we draw $\boldsymbol{\beta}$ from a stick-breaking process and then draw a series $\{\mu_m\}_{m=1}^{M}$ of i.i.d. random measures from the above process, setting $\mathbf{w} = \alpha\boldsymbol{\beta}$ for some $\alpha > 0$, then this defines a Hierarchical Dirichlet Process. If, moreover, there is one $\mu$ associated with every state $j$, then we obtain the transition prior for the HDP-HMM, where

$$p(z_t \mid z_{t-1}, \boldsymbol{\pi}) = \tilde{\pi}_{z_{t-1}z_t} \tag{7}$$

### 2.2  Promoting "Local" Transitions

In the preceding formulation, the $\theta_j$ and the $\pi_{jj'}$ are independent conditioned on the top-level measure. Our goal is to relax this assumption, in order to incorporate possible prior knowledge that certain "location" pairs, $(\theta_j, \theta_{j'})$, are more likely than others to produce large transition weights (i.e., states adjacent in time should tend to be similar). This can be accomplished by scaling the elements $\pi_{jj'}$ by a function of $(\theta_j, \theta_{j'})$ prior to normalization, or equivalently letting the Gamma distribution

2

have a proximity-dependent rate parameter. Let $\Phi : \Omega \times \Omega \to [0, \infty)$ represent a "similarity function", and define a collection of random variables $\{\phi_{jj'}\}_{j,j' \geq 1}$ according to $\phi_{jj'} = \phi(\theta_j, \theta'_j)$. We can then generalize (6) to

$$\pi_{jj'} \,|\, \boldsymbol{\beta}, \boldsymbol{\theta} \sim \mathsf{Gamma}(\alpha\beta_{j'}, \phi_{jj'}^{-1}) \qquad T_j = \sum_{j'=1}^{\infty} \pi_{jj'} \qquad \tilde{\pi}_{jj'} = \frac{\pi_{jj'}}{T_j} \tag{8}$$

so that the prior mean of $\pi_{jj'}$ is $\alpha\beta_{j'}\phi_{jj'}$. Since a similarity between one object and another should not exceed the similarity between an object and itself, and since a constant rescaling of the similarity will be absorbed in normalization, we will assume that $0 \leq \phi_{jj'} \leq 1$ for all $j$ and $j'$.

## 2.3 The HDP-HMM-LT as the Marginalization of a Markov Jump Process with "Failed" Jumps

We can gain stronger intuition, as well as simplify posterior inference, by representing the HDP-HMM-LT described in the last section as a continuous time Markov jump process where holding times have been integrated out. In particular, suppose that some of the attempts to jump from one state to another fail, and the failure probability increases as a function of the "distance" between the states.

Let $\Phi$ be defined as in the last section, and let $\boldsymbol{\beta}$, $\boldsymbol{\theta}$ and $\boldsymbol{\pi}$ be defined as in the Normalized Gamma Process representation of the ordinary HDP-HMM (so, $\pi_{jj'} \,|\, \boldsymbol{\beta}, \boldsymbol{\theta} \sim \mathsf{Gamma}(\alpha\beta_{j'}, 1)$). Now suppose that when the process is in state $j$, jumps to state $j'$ are made at rate $\pi_{jj'}$. This defines a continuous-time Markov Process where the off-diagonal elements of the transition rate matrix are the off diagonal elements of $\boldsymbol{\pi}$. In addition, self-jumps are allowed, and occur with rate $\pi_{jj}$. If we only observe the jumps and not the durations between jumps, this is an ordinary Markov chain. If we do not observe the jumps themselves, but instead an observation is generated once per jump from a distribution that depends on the state being jumped to, then we have an ordinary HMM.

We modify this process as follows. Suppose that each jump attempt from state $j$ to state $j'$ has a chance of failing, which is an increasing function of the "distance" between the states. In particular, let the success probability be $\phi_{jj'}$ (recall that we assumed above that $0 \leq \phi_{jj'} \leq 1$ for all $j, j'$). Then, the rate of successful jumps from $j$ to $j'$ is $\pi_{jj'}\phi_{jj'}$, and the corresponding rate of unsuccessful jump attempts is $\pi_{jj'}(1-\phi_{jj'})$. We denote the overall rate of successful jumps while in state $j$ overall by $T_j := \sum_{j'} \pi_{jj'}\phi_{jj'}$. Given that the process is in state $j$ at discretized time $t$ (that is, $z_t = j$), the probability that the first successful jump is to state $j'$ (that is, $z_{t+1} = j'$) is proportional to the rate of successful jump attempts to $j'$, which is $\pi_{jj'}\phi_{jj'}$. The holding time, $\tau_{t-1}$, is independent of $z_{t+1}$ and is distributed $\mathsf{Exp}(T_j)$. The total time spent in state $j$ given that it is visited $n_j$ times, is then

$$u_j \,|\, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\theta} \stackrel{ind}{\sim} \mathsf{Gamma}(n_{j\cdot}, T_j) \tag{9}$$

During this period there will be $q_{jj'}$ unsuccessful attempts to jump to state $j'$, where $q_{jj'}$ is distributed $\mathsf{Pois}(u_j\pi_{jj'}(1-\phi_{jj'}))$. Incorporating $\mathbf{u} = \{u_j\}$ as augmented data simplifies the likelihood for the transition parameters, yielding

$$L(\boldsymbol{\pi}, \boldsymbol{\phi} \,|\, \mathbf{z}, \mathbf{u}, \mathbf{Q}) = \left(\prod_{t=1}^{T} p(z_t \,|\, z_{t-1}, \boldsymbol{\pi}, \boldsymbol{\phi})\right) \prod_j p(u_j \,|\, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\phi}) \prod_{j'} p(q_{jj'} \,|\, u_j\pi_{jj'}, \phi_{jj'})$$

$$\propto \prod_j \prod_{j'} \pi_{jj'}^{n_{jj'}+q_{jj'}} \phi_{jj'}^{n_{jj'}} (1-\phi_{jj'})^{q_{jj'}} e^{-\pi_{jj'}u_j} \tag{10}$$

## 2.4 An HDP-HSMM-LT modification

We note that it is trivial to modify the HDP-HMM-LT to allow for non-Geometric duration distributions, by simply fixing the diagonal elements of $\boldsymbol{\pi}$ to be zero, allowing $D_t$ observations to be emitted $i.i.d.$ $F(\theta_{z_t})$ at jump $t$, where $D_t$ is drawn from a state-specific duration distribution, and sampling the latent state sequence using a message passing algorithm suited for HSMMs [3]. Inference for the $\boldsymbol{\phi}$ matrix is not affected, since the diagonal elements are assumed to be 1. Unlike in the original representation of the HDP-HSMM, there is no need to introduce additional auxiliary variables as a result of this modification, due to the presence of the (continuous) durations, $\mathbf{u}$, which were already needed to account for the normalization of the $\boldsymbol{\pi}$.

## 3   Inference

We develop a Gibbs sampling algorithm based on the Markov Process with Failed Jumps representation, augmenting the data with the duration variables $\mathbf{u}$, the failed jump attempt count matrix, $\mathbf{Q}$, as well as additional auxiliary variables which we will define below. In this representation the transition matrix is not represented directly, but is a function of the unscaled transition matrix $\pi$ and the similarity matrix $\phi$. The full set of variables is partitioned into blocks: $\{\gamma, \alpha, \boldsymbol{\beta}, \boldsymbol{\pi}\}$, $\{\mathbf{z}, \mathbf{u}, \mathbf{Q}, \Lambda\}$, $\{\boldsymbol{\theta}\}$, and $\{\xi\}$, where $\Lambda$ represents a set of auxiliary variables that will be introduced below, $\boldsymbol{\theta}$ represents the emission and state location parameters (which may be further factored depending on the specific choice of model), and $\xi$ represents additional parameters such as any free parameters of the similarity function, $\Phi$, and any hyperparameters of the emission distribution.

### 3.1   Sampling Transition Parameters and Hyperparameters

The joint posterior over $\gamma$, $\alpha$, $\boldsymbol{\beta}$ and $\boldsymbol{\pi}$ given the other variables will factor as

$$p(\gamma, \alpha, \boldsymbol{\beta}, \boldsymbol{\pi}) = p(\gamma)p(\alpha)p(\boldsymbol{\beta} \,|\, \gamma)p(\boldsymbol{\pi} \,|\, \alpha, \boldsymbol{\beta}) \tag{11}$$

where we have omitted the dependence on the augmented data, $\mathcal{D} = (\mathbf{z}, \mathbf{u}, \mathbf{Q}, \Lambda)$ for conciseness. We describe these four factors in reverse order.

**Sampling $\boldsymbol{\pi}$**   Having used data augmentation to simplify the likelihood for $\boldsymbol{\pi}$ to the factored conjugate form in (10), the individual $\pi_{jj'}$ are *a posteriori* independent Gamma distributed:

$$\pi_{jj'} \,|\, \alpha, \beta_{j'}, \mathcal{D} \overset{ind}{\sim} \mathsf{Gamma}(\alpha\beta_{j'} + n_{jj'} + q_{jj'}, 1 + u_j) \qquad j, j' \geq 1 \tag{12}$$

**Sampling $\boldsymbol{\beta}$**   To enable joint sampling of the latent state sequence, we employ a weak limit approximation to the HDP [3], approximating the stick-breaking process for $\boldsymbol{\beta}$ using a finite Dirichlet distribution with a finite number of components, $J$, which is larger than we expect to need. Due to the product of Gammas form, we can integrate out $\boldsymbol{\pi}$ analytically from $p(\boldsymbol{\pi}, \mathcal{D} \,|\, \boldsymbol{\beta})$, to obtain the the marginal likelihood for $\boldsymbol{\beta}$. Together, we have

$$p(\boldsymbol{\beta} \,|\, \gamma) = \frac{\Gamma(\gamma/J)^J}{\Gamma(\gamma)} \prod_j \beta_j^{\frac{\gamma}{J}-1} \qquad p(\mathcal{D} \,|\, \boldsymbol{\beta}, \alpha) \propto \prod_{j=1}^{J}(1 + u_j)^{-\alpha} \prod_{j'} \frac{\Gamma(\alpha\beta_{j'} + n_{jj'} + q_{jj'})}{\Gamma(\alpha\beta_{j'})} \tag{13}$$

where we have used the fact that the $\beta_j$ sum to 1 to pull out terms of the form $(1 + u_j)^{-\alpha\beta_{j'}}$ from the inner product in the likelihood. Following Teh et al. (2006), we can introduce auxiliary variables $\{m_{jj'}\}$, with

$$p(m_{jj'} \,|\, \beta_{j'}, \alpha, \mathcal{D}) \overset{ind}{\propto} s(n_{jj'} + q_{jj'}, m_{jj'})\alpha^{m_{jj'}} \beta_{j'}^{m_{jj'}} \tag{14}$$

for integer $m_{jj'}$ ranging between 0 and $n_{jj'} + q_{jj'}$, where $s(n, m)$ is an unsigned Stirling number of the first kind. The normalizing constant in this distribution cancels the ratio of Gamma functions in the $\boldsymbol{\beta}$ likelihood, so, letting $m_{\cdot j} = \sum_{j'} m_{j'j}$, we obtain simply

$$\boldsymbol{\beta} \,|\, \mathbf{M}, \gamma \sim \mathrm{Dirichlet}(\frac{\gamma}{J} + m_{\cdot 1}, \ldots, \frac{\gamma}{J} + m_{\cdot J}) \tag{15}$$

**Sampling Concentration Parameters**   After incorporating $\mathbf{M}$ into $\mathcal{D}$, we can integrate out $\boldsymbol{\beta}$, from the joint likelihood $p(\boldsymbol{\beta}, \mathcal{D} \,|\, \gamma, \alpha)$:

$$p(\mathcal{D} \,|\, \alpha, \gamma) \propto \alpha^{m_{\cdot\cdot}} e^{-\sum_{j''} \log(1+u_{j''})\alpha} \frac{\Gamma(\gamma)}{\Gamma(\gamma + m_{\cdot\cdot})} \prod_j \frac{\Gamma(\frac{\gamma}{J} + m_{\cdot j})}{\Gamma(\frac{\gamma}{J})} \tag{16}$$

Assume that $\alpha$ and $\gamma$ have Gamma priors. Then the update to alpha is conjugate,

$$\alpha \,|\, \mathcal{D} \sim \mathsf{Gamma}(a_\alpha + m_{\cdot\cdot}, b_\alpha + \sum_j \log(1 + u_j)), \tag{17}$$

and to simplify the likelihood for $\gamma$, we can introduce a final set of auxiliary variables, $\mathbf{r} = (r_1, \ldots, r_J)$, $r_j \in \{0, \ldots, m_{\cdot j}\}$ and $t \in (0, 1)$ with the following distributions:

$$p(r_j \,|\, m_{\cdot j}, \gamma) \propto s(m_{\cdot j}, r) \left(\frac{\gamma}{J}\right)^r \qquad p(t \,|\, m_{\cdot\cdot}\gamma) \propto t^{\gamma-1}(1 - t)^{m_{\cdot\cdot}-1}. \tag{18}$$

The normalizing constants are ratios of Gamma functions, which cancel those in (16), so that

$$\gamma \,|\, \mathcal{D} \sim \mathsf{Gamma}(a_\gamma + r_{\cdot}, b_\gamma - \log(t)) \tag{19}$$

4

## 3.2 Sampling z and the auxiliary variables

We sample the hidden state sequence, $\mathbf{z}$, jointly with the auxiliary variables, which consist of $\mathbf{u}$, $\mathbf{Q}$, $\mathbf{M}$, $\mathbf{r}$ and $t$. The joint conditional distribution of these variables is defined directly by the generative model:

$$p(\mathcal{D}) = p(\mathbf{z})p(\mathbf{u} \,|\, \mathbf{z})p(\mathbf{Q} \,|\, \mathbf{u})p(\mathbf{M} \,|\, \mathbf{z}, \mathbf{Q})p(\mathbf{r} \,|\, \mathbf{M})p(t \,|\, \mathbf{M}) \tag{20}$$

Since we are conditioning on the transition matrix, we can sample the entire sequence $\mathbf{z}$ at once with the forward-backward algorithm, as in an ordinary HMM, or its corresponding generalization if we are using the HSMM variant. Having done this, we can sample $\mathbf{u}$, $\mathbf{Q}$, $\mathbf{M}$, $\mathbf{r}$ and $t$ from their forward distributions.

## 3.3 Sampling state and emission parameters

Depending on the application, the similarities $\{\phi_{jj'}\}$ may be based directly on the emission distributions, or may be based on a separate set of variables. In the experiments described below we assume the former. For simplicity, we denote the collection of these variables by $\boldsymbol{\theta}$. We have two likelihood components:

$$p(\mathbf{z}, \mathbf{Q} \,|\, \boldsymbol{\theta}) \propto \prod_j \prod_{j'} \phi_{jj'}^{n_{jj'}} (1 - \phi_{jj'})^{q_{jj'}} \qquad p(\mathbf{Y} \,|\, \mathbf{z}, \boldsymbol{\theta}) = \prod_{t=1}^{T} f(\mathbf{y}_t; \theta_{z_t}) \tag{21}$$

where proportionality is with respect to variation in $\boldsymbol{\theta}$.

# 4 Experiments

The parameter space for the hidden states, the associated prior $H$ on $\boldsymbol{\theta}$, and the similarity function $\Phi$, is application-specific; we consider here the case where a state consists of a finite $D$-dimensional binary vector, $\eta_j$, the similarity function is a Laplacian kernel defined with respect to Hamming distance between pairs $\eta_j$ and $\eta_{j'}$ with decay parameter $\lambda$, and the emission distribution is linear-Gaussian, with $D \times K$ weight matrix $\mathbf{W}$, so that each $K$-dimensional observation is $\mathcal{N}(\mathbf{W}\eta_j, \Sigma)$. For experiments discussed here, we will assume that $\Sigma$ does not depend on $j$, but this assumption is easily relaxed if appropriate. For finite-length binary vector states, the set of possible states is finite, and so on its face it may seem that a nonparametric model is unnecessary. However, if $D$ is reasonably large, it is likely that most of the $2^D$ possible states are vanishingly unlikely (and, in fact, the number of observations may well be less than $2^D$), and so we would like a model that encourages the selection of a sparse set of states. Moreover, there could be more than one state with the same emission parameters, but with different transition dynamics. Before describing individual experiments, we describe the additional inference steps needed for these variables.

## 4.1 Additional Inference Steps

**Sampling $\eta$** We put independent Beta-Bernoulli priors on each coordinate of $\eta$. We Gibbs sample each coordinate $\eta_{jd}$ conditioned on all the others and the coordinate-wise prior means, $\{\mu_d\}$, which we sample in turn conditioned on the $\eta$s.

**Sampling $\lambda$** The Laplacian kernel $\Phi$ is defined as $\Phi(\eta_j, \eta_{j'}) = e^{-\lambda d(\eta_j, \eta_{j'})}$, where in our case $d$ is Hamming distance. The parameter $\lambda$ governs the connection between $\boldsymbol{\theta}$ and $\phi$. Writing (21) in terms of $\lambda$ and the distance matrix $\boldsymbol{\Delta}$ gives the likelihood

$$p(\mathbf{z}, \mathbf{Q} \,|\, \lambda, \boldsymbol{\theta}) \propto \prod_j \prod_{j'} e^{-\lambda \Delta_{jj'} n_{jj'}} (1 - e^{-\lambda \Delta_{jj'}})^{q_{jj'}} \tag{22}$$

We put an $\mathsf{Exp}(b_\lambda)$ prior on $\lambda$, which yields a posterior density

$$p(\lambda \,|\, \mathbf{z}, \mathbf{Q}, \boldsymbol{\theta}) \propto e^{-(b_\lambda + \sum_j \sum_{j'} \Delta_{jj'} n_{jj'})\lambda} \prod_j \prod_{j'} (1 - e^{-\lambda \Delta_{jj'}})^{q_{jj'}} \tag{23}$$

This density is log-concave, and so we use Adaptive Rejection Sampling [2] to sample from it.

5

**Sampling W and $\Sigma$**    Conditioned on the state matrix $\boldsymbol{\theta}$ and the data matrix $\mathbf{Y}$, the weight matrix $\mathbf{W}$ can be sampled as well using standard methods for Bayesian linear regression. We place a zero mean Normal prior on each element of $\mathbf{W}$ (including a row of intercept terms), resulting in a multivariate Normal posterior for each column. For the experiments reported below, we constrain $\Sigma$ to be a diagonal matrix, and place an Inverse Gamma prior on the variances, resulting in conjugate updates.

## 4.2 "Cocktail Party" Data

To evaluate the model, we created synthetic data based on a speaker diarization task, with the property that speakers are grouped into conversations, and take turns speaking within conversation. In such a task, there are naively $2^S$ possible states, where $S$ is the total number of speakers, corresponding to who is speaking when, but due to the conversational grouping, if zero or one speakers in a conversation can be speaking at any given time, the state space is constrained, with only $\prod_c (s_c + 1)$ states possible, where $s_c$ is the number of speakers in conversation $c$.

For the synthetic data, the turn sequence within conversations is generated using a Poisson HSMM with $s_c$ states, with pauses with shorter Poisson duration inserted between each "sentence". The states within conversations were then mapped to a $s_c$ length binary vector, where all zeroes corresponds to silence, and speaker $s$ speaking corresponds to a 1 in position $s$. The binary vectors were concatenated across conversations to yield latent states consisting of length $S$ binary vectors. To simulate speakers being recorded by $K$ microphones, weights from speakers to microphones were generated independently from a $U(0, 1)$ distribution, resulting in a $D \times K$ weight matrix, $\mathbf{W}$. An extra row of "background noise" parameters was added as well, also sampled from $U(0, 1)$. Independent $\mathcal{N}(0, \sigma_k^2)$ noise was added to each time step at microphone $k$.

We generated transition and emission parameters from conjugate priors to the Poisson HSMM. The data set consisted of four conversations of four speakers each, and 12 microphones, so that $D = 16$, and $K = 12$. There are therefore $2^{16} = 65536$ possible binary vector-valued states, but only $(4 + 1)^4 = 625$ can actually occur. The noise variance $\sigma_k^2$ was set to a constant of $1/10$ for all $k = 1, \ldots, K$.

We attempted to infer the states from the data using three models: (1) a binary-state Factorial HMM, in which the individual binary speaker sequences are modeled as independent a priori, (2) an ordinary HDP-HMM without local transitions, where the latent states are binary vectors, and (3) our HDP-HMM-LT model. To simplify interpretation of the results, the weight matrix was fixed to the true value (this makes the latent dimensions identifiable and makes distances between inferred and ground truth state matrices meaningful). We evaluated the models at each iteration using Hamming distance between inferred and ground truth state matrices and F1 score The results for the three models are in Figure 1. The LT model outperforms the other two on all measures on all datasets.

We also plot the number of states used by the two HDP models in Figure **??**, and the inferred decay rate $\lambda$ for the HDP-HMM-LT model. The LT model settles on a non-negligible $\lambda$ value for this data, suggesting that the local transition structure explains the data well. It also uses more components than the non-LT model, perhaps owing to the fact that the weaker transition prior of the non-LT model is more likely to explain nearby similar observations as a single persisting state, whereas the LT model places a higher probability on transitioning to a new state with a similar latent vector.

## 4.3 Synthetic Data Without Local Transitions

We also generated data from an ordinary HDP-HMM, with no local transition property, in order to investigate the performance of our model in a case where the data did not have the key property that its prior equipped it to discover. The results are in Figs. 2 and **??**. While the $\lambda$ parameter is large, the LT model has worse performance than the non-LT model; however, over iterations, the $\lambda$ parameter tends toward zero as the model learns that local transitions are not more probable, resulting in improved results.
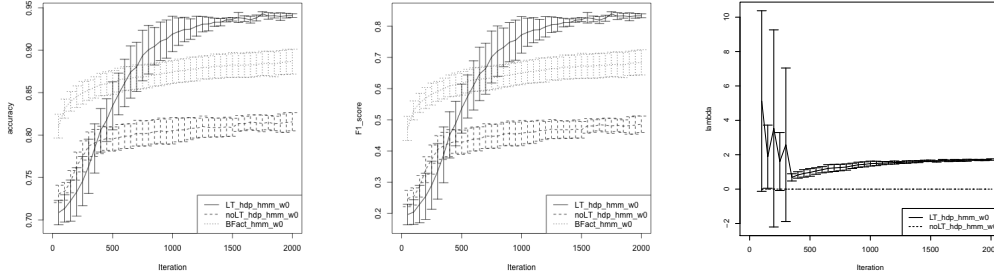
Figure 1: (a-b) Performance of the HDP-HMM-LT, standard HDP-HMM, and Binary Factorial HMM on the Cocktail Party Data. Metrics are averaged over 10 Gibbs runs on each model, with error bars representing a 99% confidence interval for the mean per iteration. The first 100 iterations are excluded as burn-in. (c) Learned similarity parameter for the LT model on the cocktail data. The first 100 iterations are excluded.
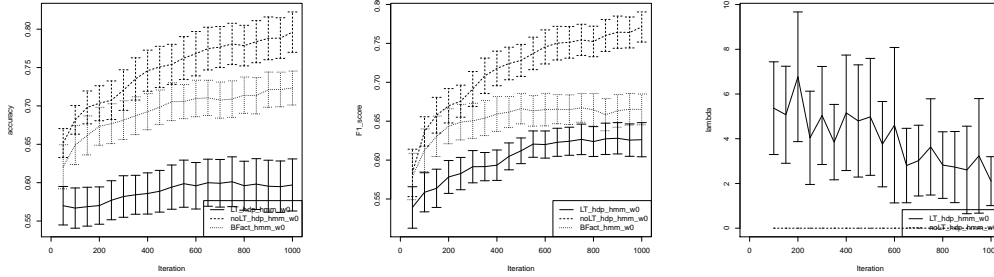


Figure 2: (a-b) Performance of the three models on data generated from an HDP-HMM without local transistions. (c) Learned similarity parameter for the LT model on the HDP-HMM data. The first 100 iterations are excluded.

## 5   Discussion

We have defined a new generative model, the Hierarchical Dirichlet Process Hidden Markov Model with Local Transitions (HDP-HMM-LT), which generalizes the HDP-HMM by allowing prior information about state space geometry to be encoded via a similarity kernel, making transitions between "nearby" pairs of states more likely *a priori*. We have derived a Gibbs sampling algorithm for this model by introducing an augmented data representation in the form of a Markov Jump Process with Failed Transitions, which also simplifies inference for the ordinary HDP-HMM. In problems with multiple dependent latent chains, the HDP-HMM-LT model combines the HDP-HMM's capacity to discover a small set of states from the large combinatorial space with the Factorial HMM's ability to encode the property that most transitions involve changes to a small number of chains at a time, outperforming both on a speaker diarization task in which speakers perform conversational groups. At the same time, despite the addition of the similarity kernel, the HDP-HMM-LT is able to learn to suppress its local transition prior when the data does not support it, performing on par with the HDP-HMM on data generated directly from the latter.

## References

[1] Emily B Fox, Erik B Sudderth, Michael I Jordan, and Alan S Willsky. An HDP-HMM for systems with state persistence. In *Proceedings of the 25th international conference on Machine learning*, pages 312–319. ACM, 2008.

[2] Walter R Gilks and Pascal Wild. Adaptive rejection sampling for gibbs sampling. *Applied Statistics*, pages 337–348, 1992.

[3] Matthew J Johnson and Alan S Willsky. Bayesian nonparametric hidden semi-markov models. *The Journal of Machine Learning Research*, 14(1):673–701, 2013.

[4] John Paisley, Chong Wang, and David M Blei. The discrete infinite logistic normal distribution. *Bayesian Analysis*, 7(4):997–1034, 2012.

[5] Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476), 2006.