

## A Doubly Hierarchical Dirichlet Process Hidden Markov Model with a Non-Ergodic Structure

Journal:	<i>Transactions on Pattern Analysis and Machine Intelligence</i>
Manuscript ID:	Draft
Manuscript Type:	Regular
Keywords:	nonparametric Bayesian models, hierarchical Dirichlet processes, hidden Markov models

SCHOLARONE™  
Manuscripts

## A Doubly Hierarchical Dirichlet Process Hidden Markov Model with a Non-Ergodic Structure

Amir H. Harati Nejad Torbati, *Student Member, IEEE*, and Joseph Picone, *Senior Member, IEEE*

**Abstract**— Nonparametric Bayesian models use a Bayesian framework to learn model complexity automatically from the data and eliminate the need for a complex model selection process. A Hierarchical Dirichlet Process Hidden Markov Model (HDPHMM) is the nonparametric Bayesian equivalent of a hidden Markov model (HMM), but is restricted to an ergodic topology that uses a Dirichlet Process Model (DPM) to achieve a mixture distribution-like model. For applications involving ordered sequences (e.g., speech recognition), it is desirable to impose a left-to-right structure on the model. In this paper, we introduce a model based on HDPHMM that: (1) shares data points between states, (2) models non-ergodic structures, and (3) models non-emitting states. The first point is particularly important because Gaussian mixture models, which support such sharing, have been very effective at modeling modalities in a signal (e.g., speaker variability). Further, sharing data points allows models to be estimated more accurately, something that is also important for an application such as speech recognition in which some mixture components occur infrequently. We demonstrate that this new model produces a 20% relative reduction in error rate for phoneme classification and an 18% relative reduction on a speech recognition task on the TIMIT Corpus.

**Index Terms**— nonparametric Bayesian models; hierarchical Dirichlet processes; hidden Markov models; speech recognition

- Corresponding Author: Amir Harati, Department of Electrical and Computer Engineering, Temple University, 1949 North 12th Street, Philadelphia, Pennsylvania, USA 19122 (Tel: 215-500-1255; Email: amir.harati@gmail.com).
- Joseph Picone is with the Department of Electrical and Computer Engineering, Temple University, Philadelphia, PA 19122. Email: joseph.picone@gmail.com.

## 1 INTRODUCTION

Hidden Markov models (HMMs) [1] are one of the most successful models for sequential data and have been applied to a wide range of applications including speech recognition. HMMs, often referred to as doubly stochastic models, are parameterized both in their structure (e.g. number of states) and emission distributions (e.g. Gaussian mixtures). Model selection methods such as the Bayesian Information Criterion (BIC) [2] are traditionally used to optimize the number of states and mixture components. However, these methods are computationally expensive and there is no consensus on an optimum criterion for selection [2].

Beal et al. [3] proposed a nonparametric Bayesian HMM with a countably infinite number of states. This model is known as an infinite HMM (iHMM) because it has an infinite number of hidden states. Teh et al. [4], [5] proposed a different formulation, HDPHMM, based on a hierarchical Dirichlet process (HDP) prior. HDPHMM is an ergodic model – a transition from an emitting state to all other states is allowed. However, in many pattern recognition applications involving temporal structure, such as speech processing, a left-to-right topology is required [7].

For example, in continuous speech recognition applications we model speech units (e.g. phonemes), which evolve in a sequential manner, using HMMs. Since we are dealing with an ordered sequence (e.g. a word is an ordered sequence of phonemes), a left-to-right model is preferred [6]. The segmentation of speech data into these units is not known in advance and therefore the training process must be able to connect these smaller models together into a larger HMM that models the entire utterance. This task can easily be achieved using left-to-right HMMs (LR-HMM). If the data has finite length, the beginning and end of a sequence is typically modeled as two additional discrete events – non-emitting initial and final states [7]. In the HDPHMM formulation, these problems are not addressed.

1  
2  
3 An HDPHMM, as well as a parametric HMM, models each emission distribution by data  
4 points mapped to that state. For example, it is common to use a Gaussian mixture model (GMM)  
5 to model the emission distributions. However, in an HDPHMM, the mixture components of these  
6 GMMs are not shared or reused. Sharing of such parameters is a critical part of most state of the  
7 art pattern recognition systems. We have introduced a model with two parallel hierarchies that  
8 enables sharing of data among different states. We refer to this model as a Doubly Hierarchical  
9 Dirichlet Process Hidden Markov Model (DHDPHMM) [8]. In this paper, we introduce a general  
10 method to add non-emitting states to both HDPHMMs and DHDPHMMs. We also develop a  
11 framework to learn non-ergodic structures from the data and present comprehensive  
12 experimental results for a standard phoneme recognition task in speech processing.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

27 The paper is organized as follows. In Section 2, we provide background on nonparametric  
28 Bayesian modeling and formally introduce the HDPHMM model. In Sections 3 and 4, we  
29 introduce the DHDPHMM model and its extensions for non-ergodic modeling and estimation of  
30 non-emitting states. In Section 5, we present results on three tasks: a pilot study on simulated  
31 data, phoneme classification and recognition on speech data. We compare DHDPHMM with  
32 both baseline and state of the art systems.  
33  
34  
35  
36  
37  
38  
39  
40

## 41 **2 BACKGROUND**

42  
43 Nonparametric Bayesian models (NPBM) have become increasingly popular in recent years  
44 because of their ability to balance model accuracy with generalization. Machine learning  
45 algorithms often have trouble dealing with previously unseen data or data sets in which the  
46 training and evaluation conditions are mismatched. Since such conditions are extremely common  
47 in applications like speech recognition, an overarching goal of this work is to improve  
48 performance when channel conditions are mismatched.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 2.1 Nonparametric Bayesian Models

A Dirichlet process (DP) [9] is a discrete distribution that consists of a countably infinite number of probability masses and defines a distribution over discrete distributions with infinite support. A DP is denoted by  $DP(\alpha, G_0)$ , and is defined as [10]:

$$G = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}, \quad \theta_k \sim G_0, \quad (1)$$

$$\beta_k = v_k \prod_{l=1}^{k-1} (1 - v_l), \quad v_k | \alpha, G_0 \sim \text{Beta}(1, \alpha). \quad (2)$$

where  $G_0$  represents the mean of the distribution [9],  $\delta_{\theta_k}$  is the unit impulse function at  $\theta_k$ ,  $\beta_k$  are weights sampled according to (2) [10], and  $\alpha$  is a concentration parameter that represents the degree of concentration around the mean ( $\alpha$  is inversely proportional to variance). The impulse functions,  $\delta_{\theta_k}$ , are often referred to as atoms.

In this representation  $\beta$  can be interpreted as a random probability measure over positive integers. The  $\beta_k$  sampled by this process, denoted by  $\beta \sim GEM(\alpha)$ , are constructed using a stick-breaking process [4]. Starting with a stick of length one, we break each stick at  $v_l$  and assign the length to  $\beta_l$ . Then we recursively break the remaining part of the stick and assign the corresponding lengths to  $\beta_k$ .

One of the main applications of a DP is to define a nonparametric prior distribution on the components of a mixture model. The resulting model is referred to as a Dirichlet Process Mixture (DPM) model and is defined as [4]:

$$\begin{aligned} \pi | \alpha &\sim GEM(\alpha) \\ z_i | \pi &\sim \text{Mult}(\pi) \\ \theta_k | G_0 &\sim G_0 \\ x_i | z_i, \{\theta_k\} &\sim F(\theta_{z_i}). \end{aligned} \quad (3)$$

In this model, the observations,  $x_i$ , are sampled from an indexed family of distributions denoted by  $F$ . If  $F$  is assumed to be Gaussian then the result is an infinite Gaussian mixture model, which is the nonparametric counterpart of a GMM [11].

An HDP extends a DPM to problems involving mixture modeling of grouped data [4] in which we desire to share components of these mixture models across groups. An HDP is defined as [4]:

$$\begin{aligned}
 G_0 &| \gamma, H \sim DP(\gamma, H) \\
 G_j &| \alpha, G_0 \sim DP(\alpha, G_0) \\
 \theta_{ji} &| G_j \sim G_j \\
 x_{ji} &| \theta_{ji} \sim F(\theta_{ji}) \quad \text{for } j \in J.
 \end{aligned} \tag{4}$$

where  $H$  provides a prior distribution for the factor  $\theta_{ji}$ ,  $\gamma$  governs the variability of  $G_0$  around  $H$  and  $\alpha$  controls the variability of  $G_j$  around  $G_0$ .  $H$ ,  $\gamma$  and  $\alpha$  are hyperparameters of the HDP. We use a DP to define a mixture model for each group and use a global DP,  $DP(\gamma, H)$ , as the common base distribution for all DPs.

## 2.2 Hierarchical Dirichlet Process Hidden Markov Model

Hidden Markov models are a class of doubly stochastic processes in which discrete state sequences are modeled as a Markov chain [1]. If we denote the state at time  $t$  with  $z_t$ , the Markovian structure can be represented by  $z_t | z_{t-1} \sim \pi_{z_{t-1}}$ , where  $\pi_{z_{t-1}}$  is the multinomial distribution that represents the transition from state  $t-1$  to state  $t$ . Observations are conditionally independent given the state of the HMM and are denoted by  $x_t | z_t \sim F(\theta_{z_t})$ . In a typical parametric HMM, the number of states is fixed so that a matrix of dimension  $N$  states by  $N$  transitions per state is used to represent the transition probabilities.

An HDPHMM is an extension of an HMM in which the number of states can be infinite. At each state  $z_t$  we can transition to an infinite number of states so the transition distribution should be drawn from a DP. However, in an HDPHMM, to obtain a chain process, we want reachable

states from one state to be shared among all states so these DPs should be linked together. In an HDPHMM each state corresponds to a group and therefore, unlike HDP in which an association of data to groups is assumed to be known a priori, we are interested in inferring this association.

A major problem with original formulation of an HDPHMM [4] is state persistence. HDPHMM has a tendency to make many redundant states and switch rapidly amongst them. Fox et al. [5] extended the definition of HDPHMM to HMMs with state persistence by introducing a sticky parameter  $\kappa$ :

$$\begin{aligned}
 \beta | \gamma &\sim GEM(\gamma) \\
 \pi_j | \alpha, \beta &\sim DP\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}\right) \\
 \psi_j | \sigma &\sim GEM(\sigma) \\
 \theta_{kj}^{**} | H, \lambda &\sim H(\lambda) \\
 z_t | z_{t-1}, \{\pi_j\}_{j=1}^{\infty} &\sim \pi_{z_{t-1}} \\
 s_t | \{\psi_j\}_{j=1}^{\infty}, z_t &\sim \psi_{z_t} \\
 x_t | \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t &\sim F(\theta_{z_t s_t}).
 \end{aligned} \tag{5}$$

The state, mixture component and observations are represented by  $z_t$ ,  $s_t$  and  $x_t$  respectively. The indices  $j$  and  $k$  are indices of the state and mixture components respectively. The base distribution,  $\beta$ , can be interpreted as the expected value of state transition distributions. The transition distribution for state  $j$  is a DP denoted by  $\pi_j$  with a concentration parameter  $\alpha$ . Another DP,  $\psi_j$ , with a concentration parameter  $\sigma$ , is used to model an infinite mixture model for each state  $z_j$ . The distribution  $H$  is the prior for the parameters  $\theta_{kj}$ .

A block sampler for HDPHMM with a multimodal emission distribution has been introduced [5] that jointly samples the state sequence  $z_{1:T}$  given the observations, model parameters and transition distribution  $\pi_j$ . A variant of the forward-backward procedure is utilized that allows us to exploit the Markovian structure of the HMM to improve the convergence speed of the

inference algorithm. However this algorithm requires approximation of the theoretically infinite distributions with a “degree  $L$  weak limit” approximation that truncates a DP into a Dirichlet distribution with  $L$  dimensions [12]:

$$GEM_L(\alpha) \triangleq Dir\left(\frac{\alpha}{L}, \dots, \frac{\alpha}{L}\right). \quad (6)$$

It should be noted that this result is different from a classical parametric Bayesian HMM since the truncated HDP priors induce a shared sparse subset of the  $L$  possible states. Interested readers can refer to [5] for more details about this algorithm.

### 3 DHDPHMM

We can extend the model in (5) to address the problem of sharable mixture components. Equation (5) defines a model with a multimodal distribution at each state. In an HDPHMM formulation these distributions are modeled using a DPM model:

$$\begin{aligned} \psi_j &| \sigma \sim GEM(\sigma) \\ s_t &| \{\psi_j\}_{j=1}^{\infty}, z_t \sim \psi_{z_t} \\ \theta_{kj}^{**} &| H, \lambda \sim H(\lambda) \\ x_t &| \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t \sim F(\theta_{z_t s_t}^{**}). \end{aligned} \quad (7)$$

Equation (7) demonstrates when the state assignment,  $z_t$ , for data point  $x_t$  is known (or sampled previously), the mixture components can be sampled from a multinomial distribution with DP priors. Equation (5) also shows that each emission distribution is modeled independent of other distributions. It has been shown previously [13] that sharing data points, if done properly, can improve the accuracy of the model.

As we have discussed in Section 2.1, HDP is the extension of a DPM to mixture modeling of grouped data. If the state assignment,  $z_t$ , is assumed to be known (or estimated) then an HDPHMM divides the data points into multiple groups. Therefore we should be able to use the



1  
2  
3 same principle and model the emission distributions with another HDP. The resulting model will  
4  
5 have two parallel hierarchies and hence is referred to as a Doubly Hierarchical Dirichlet Process  
6  
7  
8 Hidden Markov Model (DHDPHMM). Applying (4) we can write:  
9

$$\begin{aligned}
 & \xi | \tau \sim GEM(\tau) \\
 & \psi_j | \sigma, \xi \sim DP(\sigma, \xi) \\
 & \theta_{kj}^{**} | H, \lambda \sim H(\lambda) \\
 & s_t | \{\psi_j\}_{j=1}^{\infty}, z_t \sim \psi_{z_t} \\
 & x_t | \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t \sim F(\theta_{z_t, s_t}),
 \end{aligned} \tag{8}$$

10  
11  
12  
13  
14  
15  
16  
17  
18 here  $\zeta$  is the DP used as the base distribution for HDP and  $\tau$  and  $\sigma$  are hyperparameters. By  
19  
20  
21 substituting (8) in (5) we can obtain a generative model for DHDPHMM:  
22  
23  
24

$$\begin{aligned}
 & \beta | \gamma \sim GEM(\gamma) \\
 & \pi_j | \alpha, \beta \sim DP(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}) \\
 & \xi | \tau \sim GEM(\tau) \\
 & \psi_j | \sigma, \xi \sim DP(\sigma, \xi) \\
 & \theta_{kj}^{**} | H, \lambda \sim H(\lambda) \\
 & z_t | z_{t-1}, \{\pi_j\}_{j=1}^{\infty} \sim \pi_{z_{t-1}} \\
 & s_t | \{\psi_j\}_{j=1}^{\infty}, z_t \sim \psi_{z_t} \\
 & x_t | \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t \sim F(\theta_{z_t, s_t}).
 \end{aligned} \tag{9}$$

25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42 DHDPHMM pools the data points while HDPHMM divides data points between different  
43  
44 states. If we don't have enough data points in a particular state or a mixture component then the  
45  
46 distribution parameters will be estimated poorly (e.g. the mean and covariance). For example, in  
47  
48 speech recognition systems we usually use features with a dimensionality of 39 which translates  
49  
50 to  $39+(39 \times 40)/2+1=820$  free parameters per Gaussian mixture component (assuming a full  
51  
52 covariance). In an HDPHMM, with no sharing of parameters, we can easily end up with an  
53  
54 intractable number of parameters.  
55  
56  
57  
58  
59  
60

### 3.1 Inference Algorithm for DHDPHMM

An inference algorithm is required to learn the model parameters from the data. One solution to this problem is the block sampler [5] discussed in the previous section. Here we present modifications of this block sampler for inference for our DHDPHMM.

Using the “degree  $L$  weak limit” approximation to DP in (6) for HDP emissions of (8) we can write the following equations (replacing  $L'$  with  $L$ ):

$$\xi | \sigma \sim Dir\left(\frac{\sigma}{L}, \dots, \frac{\sigma}{L}\right) \quad (10)$$

$$\psi_j | \xi, \tau \sim Dir(\tau \xi_1, \dots, \tau \xi_{L'}). \quad (11)$$

Following a similar approach in [5] we write the posterior distributions for these equations as:

$$\xi | M', \tau \sim Dir\left(\frac{\tau}{L'} + M'_{\cdot 1}, \dots, \frac{\tau}{L'} + M'_{\cdot L'}\right) \quad (12)$$

$$\psi_j | \sigma, \xi, Z_{1:T}, S_{1:T} \sim Dir(\sigma \xi_1 + n'_{j1}, \dots, \sigma \xi_{L'} + n'_{jL'}) \quad (13)$$

where  $M'_{jk}$  is the number of clusters in state  $j$  with mixture component  $k$ ;  $M'_{\cdot k}$  is total number of clusters that contain mixture component  $k$ . The number of observations in state  $j$  that are assigned to component  $k$  is denoted by  $n'_{jk}$ . The posterior distribution for  $\tau$ , the hyperparameter in (12), can be written as:

$$P(\tau | n_{\cdot}, \dots, n_{J\cdot}, M'_{\cdot 1}, \dots, M'_{\cdot J\cdot}) \propto Gamma\left(a + M'_{\cdot \cdot} - \sum_{j=1}^{L'} s_j b - \sum_{j=1}^{L'} \log r_j\right) \quad (14)$$

$$P(r_j | \tau, r_{\setminus j}, s, n_{\cdot}, \dots, n_{J\cdot}, M'_{\cdot 1}, \dots, M'_{\cdot J\cdot}) \propto Beta(\tau + 1, n_{j\cdot}) \quad (15)$$

$$P(s_j | \tau, s_{\setminus j}, r, n_{\cdot}, \dots, n_{J\cdot}, M'_{\cdot 1}, \dots, M'_{\cdot J\cdot}) \propto Ber\left(\frac{n_{j\cdot}}{n_{j\cdot} + \tau}\right) \quad (16)$$

where  $r$  and  $s$  are auxiliary variables used to facilitate the inference for  $\tau$  (following the same

approach as in [5]) and  $a$  and  $b$  are hyperparameters over a Gamma distribution.

### 3.2 Scalability

The main motivation behind DHDPHMM is the ability to share mixture components and therefore data points between different states. When using the modified block sampler algorithm we only deal with  $L'$  Gaussian distributions. The HDPHMM model has  $L \times L'$  Gaussians to estimate. Since up to 95% of the inference time is spent in calculating the likelihood of data for Gaussian distributions, a reduction from  $L \times L'$  to  $L'$  reduces the computational time considerably. Also we have utilized parallel programming facilities (e.g. openMP) for the implementation of both algorithms, which makes this process feasible for large data sets. Fig. 1 provides a comparison of both algorithms for different values of  $L$  and  $L'$ . DHDPHMM's computational complexity is flat as the maximum bound on the number of states increases while the inference cost for HDPHMM grows linearly.

## 4 DHDPHMM WITH A NON-ERGODIC STRUCTURE

A non-ergodic structure for the DHDPHMM can be achieved by modifying the transition distributions. These modifications can also be applied to HDPHMM using a similar approach.

### 4.1 Left-to-Right DHDPHMM

The transition probability from state  $j$  has infinite support and can be written as:

$$\pi_j | \alpha, \beta \sim DP(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}). \quad (17)$$

From (17) we can see a transition distribution has no topological restrictions and therefore (5) and (9) define ergodic HMMs. In order to obtain a left-to-right (LR) topology we need to force the base distribution of the Dirichlet distribution in (17) to only contain atoms to the right of the current state. This means  $\beta$  should be modified so that the probability of transiting to states left of the current state (i.e. states previously visited) becomes zero. For state  $j$  we define  $V_j = \{V_{ji}\}$ :

$$V_{ji} = \begin{cases} 0, & i < j \\ 1, & i \geq j \end{cases} \quad (18)$$

where  $i$  is the index for all following states. We can then modify  $\beta$  by multiplying it with  $V_j$ :

$$\beta' = \frac{\beta \cdot V_j}{\sum_i \beta_i V_{ji}}. \quad (19)$$

In the block sampler algorithm, we have:

$$\pi_j \sim \text{Dir}(\alpha\beta'_1 + n_{j1}, \dots, \alpha\beta'_j + \kappa + n_{jj}, \dots, \alpha\beta'_L + n_{jL}), j=1, \dots, L \quad (20)$$

where  $n_{jk}$  are the number of transitions from state  $j$  to  $k$ . From (20) we can see that multiplying  $\beta$  with  $V_j$  biases  $\pi_j$  toward a left-to-right structure but there is still a positive probability to transit to the states left of  $j$ . If we leave  $\pi_j$  as in (20) the resulting model would be an LR model with possible loops. The model would be biased toward an LR structure but with the possibility of forming loops. Models with an LR structure and possible loops will be denoted as LR-L.

In order to obtain an LR model with no loops, we have to multiply  $n_{jk}$  with  $V_j$ :

$$\pi_j \sim \text{Dir}(\alpha\beta'_1 + V_{j1}n_{j1}, \dots, \alpha\beta'_j + \kappa + V_{jj}n_{jj}, \dots, \alpha\beta'_L + V_{jL}n_{jL}), j=1, \dots, L. \quad (21)$$

$V_j$  and  $\beta'$  are calculated from (18) and (19) respectively. This model always finds transitions to the right of state  $j$  and is referred to as an LR model.

Sometimes it is useful to have LR models that allow restricted loops to the first state. For example, when dealing with long sequences, a sequence might have a local left to right structure but needs a reset at some point in time. To modify  $\beta$  to obtain an LR model with a loop to the first state (LR-LF) we can write:

$$V_{ji} = \begin{cases} 0, & 0 < i < j \\ 1, & i \geq j, i=0 \end{cases} \quad (22)$$

$\beta'$  can be calculated from (19) and  $\pi_j$  should be sampled from (21).

The LR models described above allow for skip transitions that mean the model learns parallel

paths that correspond to different modalities in the training data. Sometimes more restrictions on the structure might be required. One such example is a strictly left to right structure (LR-S):

$$V_{ji} = \begin{cases} 0, & i \neq j+1 \\ 1, & i = j+1 \end{cases} \quad (23)$$

## 4.2 Initial and Final Non-Emitting States

In many applications, such as speech recognition, an LR-HMM begins from and ends with non-emitting states. These states are required to model the beginning and end of finite duration sequences. Adding a non-emitting initial state is straightforward: the probability of transition into the initial state is 1 and the probability distribution of a transition from this state is equal to  $\pi_{init}$  which is the initial probability distribution for an HMM without non-emitting states. However, adding a final non-emitting state is more complicated. In the following sections we will discuss two approaches that solve this problem.

### 4.2.1 Maximum Likelihood Estimation

Consider state  $z_i$  depicted in Fig. 2. The outgoing probabilities for any state can be classified into three categories: (1) a self-transition ( $P_1$ ), (2) a transition to all other states ( $P_2$ ), and (3) a transition to a final non-emitting state ( $P_3$ ). These probabilities must sum to 1:  $P_1+P_2+P_3=1$ . Suppose that we obtain  $P_2$  from the inference algorithm. We will need to reestimate  $P_1$  and  $P_3$  from the data. This problem is, in fact, equivalent to the problem of tossing a coin until we obtain the first tails. Each head is equal to a self-transition and the first tails triggers a transition to the final state. This can be modeled using a geometric distribution [14]:

$$P(x = k) = (1 - \rho)^{k-1} \rho. \quad (24)$$

Equation (24) shows the probability of  $K-1$  heads before the first tail. In this equation  $1-\rho$  is the probability of heads (success). We also have:

$$\frac{P_1}{1-P_2} = 1-\rho, \quad \frac{P_3}{1-P_2} = \rho. \quad (25)$$

Suppose we have a total of  $N$  examples but for a subset of these,  $M_i$ , the state  $z_i$  is the last state of the model ( $S_M$ ). It can be shown [14] that the maximum likelihood estimation is obtained by:

$$\hat{\rho}_i = \frac{M_i}{\sum_{j \in S_M} k_j} \quad (26)$$

where  $k_i$  are the number of self-transitions for state  $i$ . Notice that if  $z_i$  is never the last state, then  $M_i = 0$  and  $P_3 = 0$ .

#### 4.2.2 Bayesian Estimation

Another approach to estimate transitions to a final non-emitting state,  $\rho_i$ , is to use a Bayesian framework. Since a beta distribution is the conjugate distribution for a geometric distribution, we can use a beta distribution with hyperparameters  $(a, b)$  as the prior and obtain a posterior as [15], [16]:

$$\rho_i \sim \text{Beta} \left( a + M_i, b + \sum_{j \in S_M} (k_j - 1) \right) \quad (27)$$

where  $M_i$  and  $S_M$  are the number of times which state  $z_i$  was the last state and set of all such states respectively. Hyperparameters  $(a, b)$  can also be estimated using a Gibbs sampler if required [17]. If we use (27) to estimate  $\rho_i$  we need to modify (20) to impose the constraint that the sum of the transition probabilities add to one. This is a relatively simple modification based on the stick-breaking interpretation of a DP in (2). This modification is equal to assigning  $\rho_i$  to the first break of the stick and then breaking the remaining  $1-\rho_i$  portion as before.

#### 4.3 An Integrated Model

The final definition for DHDPHMM model with a non-ergodic structure is given by:

$$\begin{aligned}
 \beta | \gamma &\sim GEM(\gamma), \beta' = \frac{V_j \cdot \beta}{\sum_i V_{ji} \beta_i} \\
 \pi_j | \alpha, \beta' &\sim DP(\alpha + \kappa, \frac{\alpha \beta' + \kappa \delta_j}{\alpha + \kappa}) \\
 \xi | \tau &\sim GEM(\tau) \\
 \psi_j | \sigma, \xi &\sim DP(\sigma, \xi) \\
 \theta_{kj}^{**} | H, \lambda &\sim H(\lambda) \\
 z_t | z_{t-1}, \{\pi_j\}_{j=1}^{\infty} &\sim \pi_{z_{t-1}} \\
 s_t | \{\psi_j\}_{j=1}^{\infty}, z_t &\sim \psi_{z_t} \\
 x_t | \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t &\sim F(\theta_{z_t, s_t}).
 \end{aligned} \tag{28}$$

In this definition  $V_i$  should be replaced with the proper definition from previous section based on the type of structure we want. For example if we want an LR model then  $V_i$  should be sampled from (18). Also note that by setting  $V_i$  to one we obtain the ergodic DHDPHMM in (9). A graphical representation is shown in Fig. 3-b. The HDPHMM [5] is also displayed in Fig. 3-a for comparison.

We have not incorporated modeling of non-emitting states discussed above in (28). If we choose to use a maximum likelihood approach for estimating the non-emitting states then no change to this model is required (e.g. we can estimate these non-emitting states after estimating other parameters). However, if we choose to use the Bayesian approach then we have to replace the sampling of  $\pi_j$  in (28) with:

$$\begin{aligned}
 \bar{w}, \bar{\chi} &\sim \text{Modified stick-breaking}(\alpha, \kappa) \\
 \pi_j | \bar{w}, \bar{\chi} &\sim \sum_k w_k \delta_{\chi_k}
 \end{aligned} \tag{29}$$

$$\text{Modified stick-breaking}(\alpha, \kappa) = \begin{cases} \text{for } i = \{1, 2, \dots\}: \\ v_i | \alpha, \kappa \sim \text{Beta}(1, \alpha + \kappa) \\ w_i | v_i, \rho_j = v_i(1 - \rho_j) \prod_{l=1}^{k-1} (1 - v_l) \\ \chi_i | \alpha, \beta', \kappa \sim \sum_k \frac{\alpha \beta'_k + \kappa \delta_{kj}}{\alpha + \kappa} \delta_k \end{cases} \tag{30}$$

1  
2  
3  
4 where we have replaced DP with the modified stick-breaking process described above.

## 5 EXPERIMENTS

6  
7  
8  
9 In this section we provide some experimental results which compare DHDPHMM with  
10 HDPHMM, HMM and several other state of the art models. The experiments begin with artificial  
11 data and then proceed to standard phoneme classification and recognition tasks.  
12  
13  
14

### 15 5.1 HMM-Generated Data

16  
17 To demonstrate the basic efficacy of the model, we generated data from a 4-state left to right  
18 HMM. The emission distribution for each state is a GMM with a maximum of three components,  
19 each consisting of a two-dimensional normal distribution. Three synthetic data sequences  
20 totaling 1900 observations were generated for training. Three configurations have been studied:  
21 (1) an ergodic HDPHMM, (2) an LR HDPHMM and (3) an LR DHDPHMM. A Normal-inverse-  
22 Wishart distribution (NIW) prior is used for the mean and covariance. The truncation levels are  
23 set to 10 for both the number of states and the number of mixture components.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34

35 Fig. 4-a shows the average likelihood for different models for held-out data by averaging five  
36 independent chains. Fig. 4-b compares the trained model to the reference structure. The LR  
37 DHDPHMM discovers the correct structure while the ergodic HDPHMM finds a more simplified  
38 HMM because LR DHDPHMM constrains the search space to left to right topologies while  
39 HDPHMM has a less constrained search space. Further, we can see that DHDPHMM has a  
40 higher overall likelihood. While LR HDPHMM can find the structure close to the correct one, its  
41 likelihood is slightly lower than the ergodic HDPHMM. However, LR DHDPHMM produces a  
42 15% (relative) improvement in likelihoods compared to the ergodic model. It is also interesting  
43 to note that the likelihoods of models discovered by all the nonparametric Bayesian algorithms  
44 are superior to the likelihood of the reference model itself.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## 5.2 Phoneme Classification on the TIMIT Corpus

The TIMIT Corpus [18] is one of the most cited evaluation data sets used to compare new speech recognition algorithms. The data is segmented manually into phonemes and therefore is a natural choice to evaluate phoneme classification algorithms. TIMIT contains 630 speakers from eight main dialects of American English. There are a total of 6,300 utterances where 3,990 are used in the training set and 192 utterances are used for the “core” evaluation subset (another 400 used as development set). We followed the standard practice of building models for 48 phonemes and then map them into 39 phonemes [20].

A standard 39-dimensional MFCC feature vector was used (12 Mel-frequency Cepstral Coefficients plus energy and their first and second derivatives) to convert speech data into feature streams. Cepstral mean subtraction [19] was also used.

### 5.2.1 A Comparison to HDPHMM

In Table 1 we compare the performance of DHDPHMM to HDPHMM. We provide error rates for both the development and core subsets. In this table we have compared an LR model with two other models: a strictly LR topology and an ergodic model. As this table shows DHDPHMM is consistently better than their HDPHMM counterparts. Further, it can be seen that LR models perform better than ergodic models (as expected) while strictly LR models perform more poorly. This is due to the fact that a strictly LR model constrains the best path to one path while the other LR models learn many parallel paths. From the last column of this table we can see LR DHDPHMM finds 3888 Gaussians for all 48 phonemes while two different LR HDPHMM models find 4628 and 7281 Gaussians for all phonemes respectively. These numbers show DHDPHMM can learn a less complex model that can explain the data better than a more complex model learned by HDPHMM. This is an important property that validates the basic

1  
2  
3 philosophy of the NPBM and also follows Occam's Razor [20].  
4

5 Fig. 5 shows the structures for phonemes /aa/ and /sh/ discovered by DHDPHMM. It is clear  
6 that the model structure evolves with amount of data points, validating another characteristic of  
7 the NPBM. It is also important to note that the structure learned for each phoneme is unique and  
8 reflects underlying differences between phonemes. Finally, note that the proposed model learns  
9 multiple parallel left-to-right paths. This is shown in Fig. 5-b where S1-S2, S1-S3 and S1-S4  
10 depict three parallel models.  
11  
12  
13  
14  
15  
16  
17  
18  
19

20 Fig. 6 show the confusion matrix for the most confusable pairs of this classification task. The  
21 general confusion matrix follows the same trend but because it is too large it has not been shown  
22 in this paper. From this confusion matrix we can see that most errors occur, as expected, between  
23 acoustically similar phonemes. In fact, if we use 5 broad phonetic classes instead of using 39  
24 phoneme classes, the classification error rate drops to 4.8%.  
25  
26  
27  
28  
29  
30  
31

### 32 **5.2.2 A Comparison to Other Representative Systems**

33 Table 2 shows a full comparison between DHDPHMM and both baseline and state of the art  
34 systems. The first three rows of this table show three-state LR HMMs trained using maximum  
35 likelihood (ML) estimation. HMM with 40 Gaussians per state performs better than other two  
36 and has an error rate of 26.1% on the core subset. Our LR DHDPHMM model has error rate of  
37 21.4% on the same subset of data (a 20% relative improvement). It should be noted that the  
38 number of Gaussians used by this HMM system is 5760 (set a priori) while our LR DHDPHMM  
39 uses only 3888 Gaussians. Fig. 7 shows the error rate vs. the amount of training data for both  
40 HMM and DHDPHMM systems. As we can see DHDPHMM is always better than the HMM  
41 model. For example, when trained only using 40% of data DHDPHMM performs better than an  
42 HMM using the entire data set. Also it is evident that HMM performance does not improve  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 significantly when we train it with more than 60% of the data (error rates for 60% and 100% are  
4 very close) while DHDPHMM improves with more data.  
5  
6

7  
8 Fig. 8 shows the number of Gaussians discovered by DHDPHMM versus the amount of data.  
9  
10 The model evolves into a more complex model as it is exposed to more data. This growth in  
11 complexity is not linear (e.g. number of Gaussians grows 33% when the amount of data  
12 increases 5 times) which is consistent with the DP prior constraints. If we want to change this  
13 behavior we would have to use other type of priors.  
14  
15  
16  
17  
18

19  
20 The fourth row of Table 2 shows the error rate for an HMM trained using a discriminative  
21 objective function (e.g. MMI). We can see discriminative training reduces the error rate.  
22  
23 However, the model still produces a larger error rate relative to our ML trained DHDPHMM.  
24  
25 This suggests that we can further improve DHDPHMM if we use discriminative training  
26 techniques. Several other state of the art systems are shown that have error rates comparable to  
27 our model. Data-driven HMMs [24], unlike DHDPHMM, models the context implicitly, which  
28 seems to be one of the main reasons that it performs so well. We expect to obtain better results if  
29 we also use context dependent (CD) models instead of context independent (CI) models.  
30  
31  
32  
33  
34  
35  
36  
37

### 38 **5.3 Supervised Phoneme Recognition**

39  
40  
41 Speech recognition systems usually use a semi-supervised method to train acoustic models. By  
42 semi-supervised we mean the exact boundaries between phonemes are not given but instead the  
43 transcription only consists of a sequence of phones in an utterance. It is has been shown that this  
44 semi-supervised method actually works better than a completely supervised method [24].  
45  
46 However, in this section we use a completely supervised method to evaluate DHDPHMM  
47 models for a phoneme recognition task. As in the previous section, we have trained  
48 DHDPHMMs only using maximum likelihood and with no context information.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 In the phoneme recognition problem, unlike phoneme classification, the boundaries between  
4 subsequent phonemes are not known (during the recognition phase) and should be estimated  
5 along with phoneme labels. During recognition we have to decide if a given frame belongs to the  
6 current group of phonemes under consideration or we have to initiate a new phoneme hypothesis.  
7 This decision is made by considering both the likelihood measurements and the language model  
8 probabilities. All systems compared in this section use bigram language models. However, the  
9 training procedure and optimization of each language model is different and has some effect on  
10 the reported error rates.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

22 In the following we define *% Correct* and *% Error* as follows [19]:  
23

$$24 \quad \% \text{ Correct} = \frac{N - S - D}{N} \quad (31)$$

$$25 \quad \% \text{ Error} = \frac{S + D + I}{N} \quad (32)$$

26  
27  
28 where  $N$  is the total number of labels in the reference transcriptions,  $S$  is the number of  
29 substitution errors,  $D$  is the number of deletion errors and  $I$  is the number of insertion errors.  
30  
31

32 Table 3 presents results for several state of the art models. As we can see, systems can be  
33 divided into two groups based on their training method (discriminative or not) and context  
34 modeling. The first two rows of this table show two similar HMM based systems with and  
35 without contextual information. We can see the error rate drops from 35.4% to 26.2% when we  
36 use a system with contextual modeling. We can also see DHDPHMM works much better than a  
37 comparable CI HMM model (the error rate drops from 35.4% for HMM to 28.6% for  
38 DHDPHMM).  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51

52 The third and fourth rows show two context-dependent HMM models. DHDPHMM performs  
53 slightly better than the CD model in row three (CD HMM 2) but slightly worse than CD model  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 of row four (CD HMM 3). We expect to obtain much better results if we use context dependent  
4  
5 models. Our model also performs better than a discriminatively trained context-independent  
6  
7 HMM. By comparing DHDPHMM with other systems presented in Table 3 we can see  
8  
9 DHDPHMM is among the best models for context-independent systems but is not as good as  
10  
11 state of the art context-dependent models.  
12  
13

## 14 15 **6 CONCLUSIONS**

16  
17 In this paper we introduced a DHDPHMM that is an extension of HDPHMM which  
18  
19 incorporates a parallel hierarchy to share data between states. We have also introduced methods  
20  
21 to model non-ergodic structures. We demonstrated through experimentation that LR  
22  
23 DHDPHMM outperforms both HDPHMM and its parametric HMM counterparts. We have also  
24  
25 shown that despite the fact that we have only used ML training for DHDPHMM performance is  
26  
27 comparable to discriminatively trained models. Further, DHDPHMM provides the best  
28  
29 performance among context-independent models.  
30  
31  
32

33  
34 Future research will focus on incorporating semi-supervised training and context modeling.  
35  
36 We have also shown that complexity grows very slowly with the data size because of the DP  
37  
38 properties (only 33% more Gaussians were used after increasing the size of the data five times).  
39  
40 Therefore it makes sense to explore other types of prior distributions to investigate how it can  
41  
42 affect the estimated complexity and overall performance. Another possible direction is to replace  
43  
44 HDP emissions with more general hierarchical structures such as a Dependent Dirichlet Process  
45  
46 [31] or an Analysis of Density (AnDe) model [32]. It has been shown that the AnDe model is the  
47  
48 appropriate model for problems involving sharing among multiple sets of density estimators [4],  
49  
50 [20].  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## ACKNOWLEDGEMENTS

The authors wish to thank Professor Marc Sobel for many valuable discussions on these topics.

This research was supported in part by the National Science Foundation through Major Research Instrumentation Grant No. CNS-09-58854.

## REFERENCES

1. L. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257-286, 1989.
2. J. B. Kadane and N. A. Lazar, "Methods and Criteria for Model Selection," *Journal of the ASA*, vol. 99, no. 465, pp. 279-290, 2004.
3. M. Beal, Z. Ghahramani, and C. E. Rasmussen, "The Infinite Hidden Markov Model," *Proceedings of NIPS*, 2002, pp. 577-584.
4. Y. Teh, M. Jordan, M. Beal, and D. Blei, "Hierarchical Dirichlet Processes," *Journal of the ASA*, vol. 101, no. 47, pp. 1566-1581, 2006.
5. E. Fox, E. Sudderth, M. Jordan, and A. Willsky, "A Sticky HDP-HMM with Application to Speaker Diarization.," *The Annals of Applied Statistics*, vol. 5, no. 2A, pp. 1020-1056, 2011.
6. B.-H. Juang and L. Rabiner, "Hidden Markov Models for Speech Recognition," *Technometrics*, vol. 33, no. 3, pp. 251-272, 1991.
7. G. A. Fink, "Configuration of Hidden Markov Models From Theory to Applications," *Markov Models for Pattern Recognition*, Springer Berlin Heidelberg, 2008, pp. 127-136.
8. A. Harati Nejad Torbati, J. Picone, and M. Sobel, "A Left-to-Right HDP-HMM with HDPM Emissions," *Proceedings of the CISS*, 2014, pp. 1-6.
9. Y.-W. Teh, "Dirichlet process," *Encyclopedia of Machine Learning*, Springer, 2010, pp. 280-287.
10. J. Sethuraman, "A constructive definition of Dirichlet priors," *Statistica Sinica*, vol. 4, no. 2, pp. 639-650, 1994.
11. C. E. Rasmussen, "The Infinite Gaussian Mixture Model," *Proceedings of NIPS*, 2000, pp. 554-560.
12. H. Ishwaran and M. Zarepour, "Exact and approximate sum representations for the Dirichlet process.," *Canadian Journal of Statistics*, vol. 30, no. 2, pp. 269-283, 2002.

13. S. Young and P. C. Woodland, "State clustering in HMM-based continuous speech recognition," *Computer Speech & Language*, vol. 8, no. 4, pp. 369-383, 1994.
14. J. Pitman, *Probability*. New York, New York, USA: Springer-Verlag, 1993, pp. 480-498.
15. A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin, *Bayesian Data Analysis*, 2nd ed. Chapman & Hall, 2004.
16. P. Diaconis, K. Khare, and L. Saloff-Coste, "Gibbs Sampling, Conjugate Priors and Coupling," *Sankhya A*, vol. 72, no. 1, pp. 136-69, 2010.
17. F. A. Quintana and W. Tam, "Bayesian Estimation of Beta-binomial Models by Simulating Posterior Densities," *Journal of the Chilean Statistical Society*, vol. 13, no. 1-2, pp. 43-56, 1996.
18. J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallet, N. Dahlgren, and V. Zue, "TIMIT Acoustic-Phonetic Continuous Speech Corpus," The Linguistic Data Consortium Catalog, 1993. [Online]. Available: <http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC93S1>
19. S. Young, G. Evermann, M. Gales, T. Hain, D. Kershaw, X. Liu, G. Moore, J. Odell, D. Ollagson, D. Povey, V. Valtchev, and P. Woodland, *The HTK Book*, Cambridge, UK, 2006.
20. C. E. Rasmussen and Z. Ghahramani, "Occam's Razor," *Proceedings of NIPS*, 2001, pp. 294-300.
21. A. Gunawardana, M. Mahajan, A. Acero, and J. C. Platt, "Hidden Conditional Random Fields for Phone Classification," *Proceedings of INTERSPEECH*, 2005, pp. 1117-1120.
22. F. Sha and L. K. Saul, "Large Margin Gaussian Mixture Modeling for Phonetic Classification and Recognition," *Proceedings of ICASSP*, 2006, pp. 265-268.
23. P. Clarkson and P. J. Moreno, "On the use of support vector machines for phonetic classification," *Proceedings of ICASSP*, 1999, pp. 585-588.
24. S. Petrov, A. Pauls, and D. Klein, "Learning Structured Models for Phone Recognition," *Proceedings of EMNLP-CoNLL*, 2007, pp. 897-905.
25. K.-F. Lee and H.-W. Hon, "Speaker-independent phone recognition using hidden Markov models," *IEEE Transactions on ASSP*, vol. 37, no. 11, pp. 1641-1648, 1989.
26. L. Lamel and J.-L. Gauvain, "High Performance Speaker-Independent Phone Recognition Using CDHMM," *Proceedings of EUROSPEECH*, 1993, pp. 121-124.
27. S. Kapadia, V. Valtchev, and S. Young, "MMI training for continuous phoneme recognition on the TIMIT database," *Proceedings of ICASSP*, 1993, pp. 491-494.

- 1
- 2
- 3
- 4 28. A.K. Halberstadt and J. R. Glass, "Heterogeneous acoustic measurements and multiple
- 5 classifiers for speech recognition," *Proceedings of ICSLP*, 1998, pp. 995-998.
- 6
- 7 29. J. Morris and E. Fosler-Lussier, "Conditional Random Fields for Integrating Local
- 8 Discriminative Classifiers," *IEEE Transactions on ASSP*, vol. 16, no. 3, pp. 617-628, 2008.
- 9
- 10 30. D. Palaz, R. Collobert, and M. Magimai-Doss, "End-to-end Phoneme Sequence Recognition
- 11 using Convolutional Neural Networks," *Proceedings of the NIPS Deep Learning Workshop*,
- 12 2013, pp. 1-8.
- 13
- 14 31. S. N. MacEachern, "Dependent Nonparametric Processes," in *ASA Proceedings of the*
- 15 *Section on Bayesian Statistical Science*, 1999, pp. 50-55.
- 16
- 17
- 18 32. G. Tomlinson and M. Escobar, "Analysis of Densities," University of Toronto, Toronto,
- 19 Canada, 1999.
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60



## List of Figures:

Fig. 1. DHDPHMM improves scalability relative to HDPHMM

Fig. 2. Outgoing probabilities for state  $z_i$

Fig. 3. Comparison of models: (a) ergodic HDPHMM [5] (b) DHDPHMM

Fig. 4. Comparison of (a) log-likelihoods of the proposed models to an ergodic model, and (b) the corresponding model structures

Fig. 5. An automatically derived model structure for a left-to-right DHDPHMM model (without the first and last non-emitting states) for (a) /aa/ with 175 examples (b) /aa/ with 2,256 examples (c) /sh/ with 100 examples and (d) /sh/ with 1,317 examples

Fig. 6. Confusion matrix for phoneme classification for the most confusable pairs

Fig. 7. Error rate vs. amount of training data for LR DHDPHMM and LR HMM

Fig. 8. Number of discovered Gaussians vs. amount of training data

## List of Tables:

TABLE 1 Comparison of LR DHDPHMM with HDPHMM

TABLE 2 Comparison of LR DHDPHMM to other algorithms

TABLE 3 Comparison of phoneme recognition performance

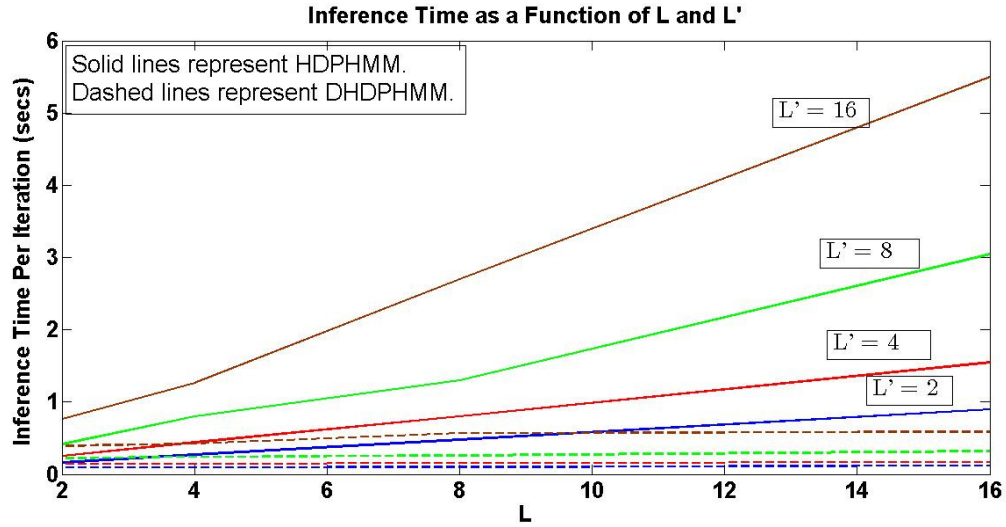


Fig. 1. DHDPHMM improves scalability relative to HDPHMM

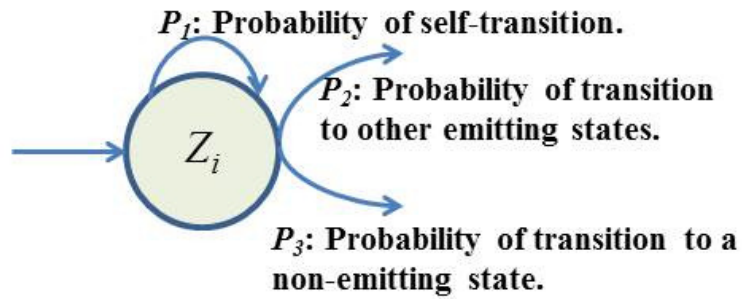
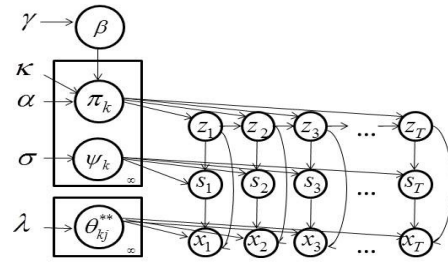
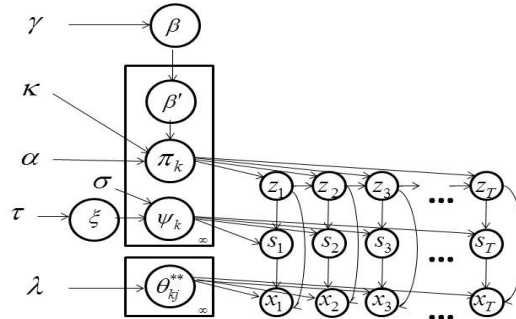


Fig. 2. Outgoing probabilities for state  $z_i$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



(a)



(b)

Fig. 3. Comparison of models: (a) ergodic HDPHMM [5] (b) DHDPHMM

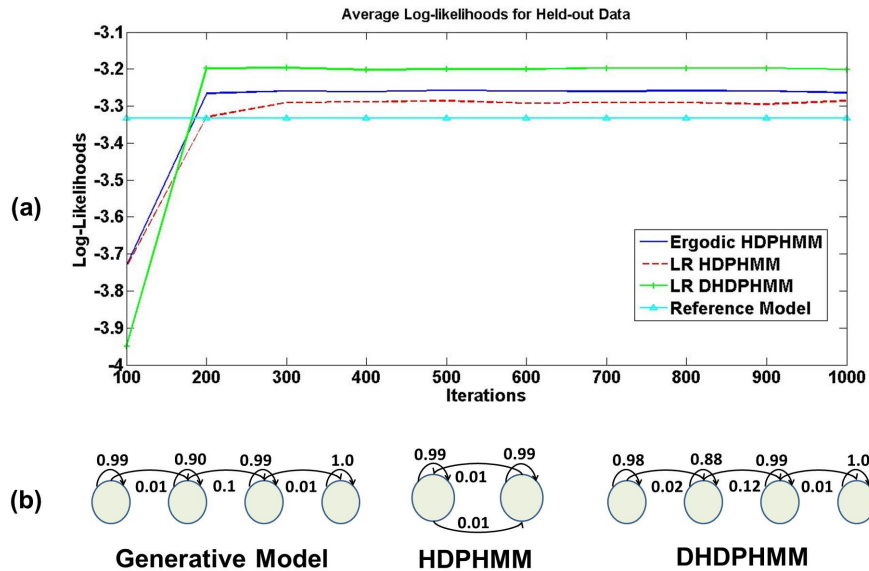


Fig. 4. Comparison of (a) log-likelihoods of the proposed models to an ergodic model, and (b) the corresponding model structures

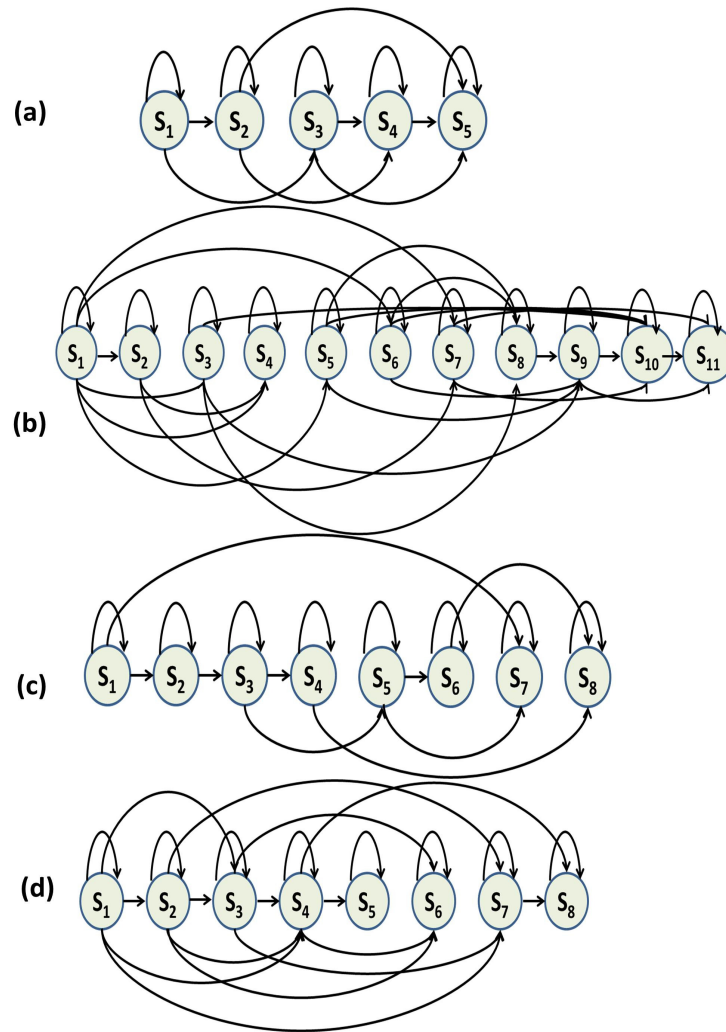


Fig. 5. An automatically derived model structure for a left-to-right DHDPHMM model (without the first and last non-emitting states) for (a) /aa/ with 175 examples (b) /aa/ with 2,256 examples (c) /sh/ with 100 examples and (d) /sh/ with 1,317 examples

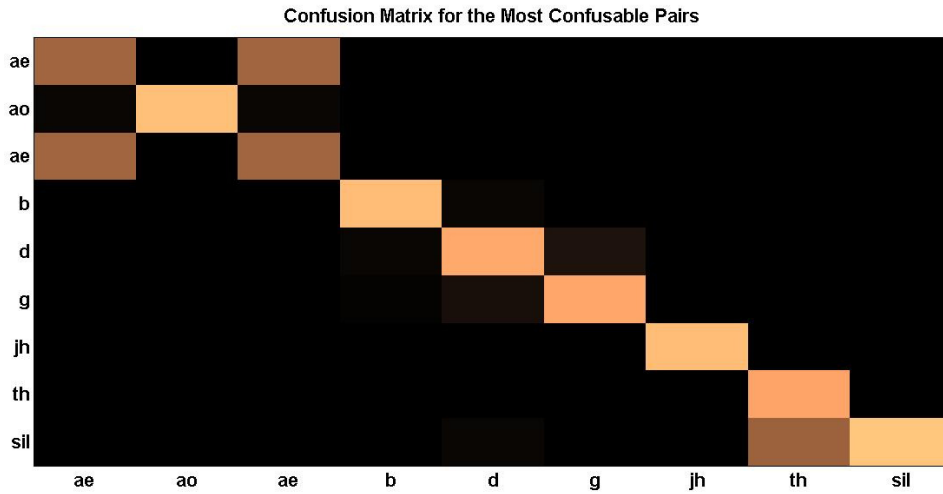


Fig. 6. Confusion matrix for phoneme classification for the most confusable pairs



Fig. 7. Error rate vs. amount of training data for LR DHDHMM and LR HMM

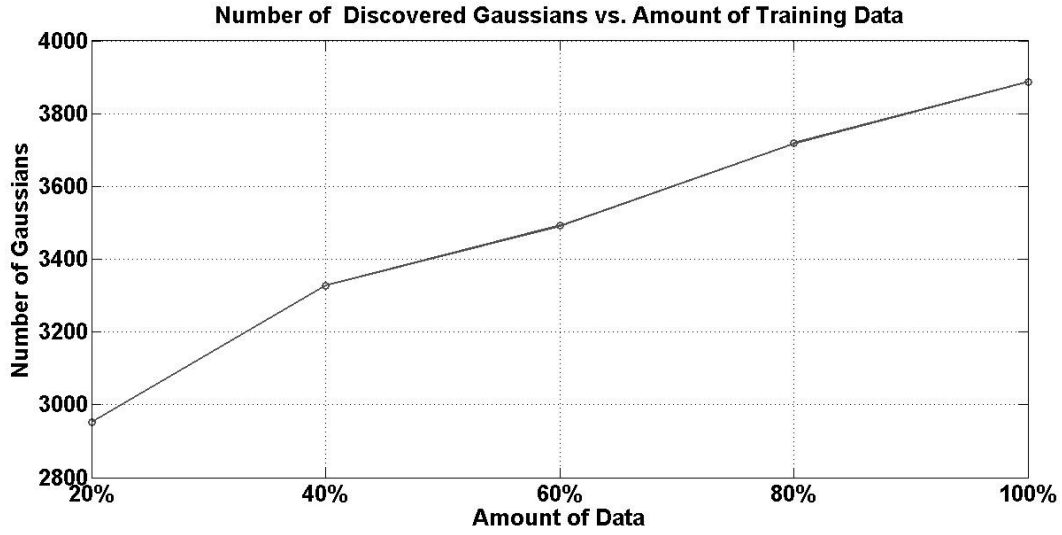


Fig. 8. Number of discovered Gaussians vs. amount of training data

TABLE 1  
COMPARISON OF LR DHDPHMM WITH HDPHMM

Model	Dev Set (% Error)	Core Set (% Error)	No. Gauss.
LR HDPHMM 1	23.5%	24.4%	4628
LR HDPHMM 2	23.8%	25.1%	7281
Ergodic DHDPHMM	24.0%	25.4%	2704
Strictly LR DHDPHMM	39.0%	38.4%	2550
<b>LR DHDPHMM</b>	<b>20.5%</b>	<b>21.4%</b>	<b>3888</b>

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

TABLE 2  
COMPARISON OF PHONEME RECOGNITION PERFORMANCE

Model	Discrim. Training	Dev Set (% Error)	Core Set (% Error)
HMM (10 Gauss.)	No	28.4%	28.7%
HMM (20 Gauss.)	No	26.1%	27.3%
HMM (40 Gauss.)	No	25.0%	26.1%
HMM/MMI (20 Gauss.) [20]	Yes	23.2%	24.6%
HCRF/SGD [20]	Yes	20.3%	21.7%
Large Margin GMMs [22]	Yes	–	21.1%
GMMs/Full Cov. [22]	No	–	26.0%
SVM [23]	Yes	–	22.4%
Data-driven HMM [24]	No	–	21.4%
<b>LR DHDPHMM</b>	<b>No</b>	<b>20.5%</b>	<b>21.4%</b>

TABLE 3  
COMPARISON OF PHONEME RECOGNITION PERFORMANCE

Model	Discrim. Training	Context Modeling	% Error	% Correct	Subset
CI-HMM [25]	No	No	35.9%	–	TID7
CD-HMM 1[25]	No	Yes	26.2%	–	TID7
CD-HMM 2[26]	No	Yes	30.9%	–	Core
CD-HMM 3[13]	No	Yes	27.7%	–	Core
HMM MMI 1 [27]	Yes	No	32.5%	73.5%	Random
HMM MMI 2 / Full Cov. [27]	Yes	No	30.3%	74.4%	Random
Heterogeneous Class. [28]	Yes	Yes	24.4%	–	Core
Data-driven HMM [24]	N/A	Yes	26.4%	–	Core
Large Margin GMM [22]	Yes	No	30.1%	–	Core
CRF [29]	Yes	No	29.9%	73.2%	All
Tandem HMM [29]	Yes	Yes	30.6%	75.6%	All
CNN/CRF [30]	Yes	No	29.9%	–	Core
<b>LR DHDPHMM</b>	<b>No</b>	<b>No</b>	<b>29.7%</b>	<b>74.1%</b>	<b>Core</b>
<b>LR DHDPHMM</b>	<b>No</b>	<b>No</b>	<b>28.6%</b>	<b>75.1%</b>	<b>Dev</b>
<b>LR DHDPHMM</b>	<b>No</b>	<b>No</b>	<b>29.2%</b>	<b>74.7%</b>	<b>All</b>

1  
2  
3 A preliminary version of this work has been published here:

4 Harati Nejad Torbati, A. H., Picone, J., & Sobel, M. (2014). A Left-to-Right HDP-HMM with HDPM  
5 Emissions. In *Proceedings of the Conference on Information Sciences and Systems* (pp. 1–6). Princeton, New  
6 Jersey, USA.  
7

8 This paper presented the initial idea, but did not include many important mathematical details or a  
9 comprehensive set of experiments.  
10

11 Major Theoretical Differences:

- 12 1. The model introduced in the conference paper was restricted to only a special case of the more  
13 general model introduced in this paper. Specifically, in this paper we have introduced a general  
14 model named DHDPHMM and its inference algorithm while in the conference paper we only  
15 introduced a special case of the model with no details about the inference algorithm.  
16
- 17 2. In this paper, we provide important theoretical derivations and implementation details regarding  
18 DHDPHMM. We discuss its differences relative to other models (e.g. HDPHMM). These details  
19 are not in the conference paper.  
20
- 21 3. In this paper, we introduce a general framework for non-Ergodic structures while in the conference  
22 paper we only derived a left to right structure.  
23
- 24 4. In this paper, we have included an inference algorithm for a Bayesian approach of adding non-  
25 emitting states while in conference paper these details were missing.  
26

27  
28 Major Experimental Differences:

- 29 1. In this paper we have a very extensive experimental section while in the conference paper the  
30 experimental section was very brief. The only commonality is section of simulated data.  
31
- 32 2. Phoneme classification experiments are much more comprehensive in this paper. In the conference  
33 paper they were limited to left to right HDPHMMs and HMMs. In this paper we have added results  
34 for DHDPHMM (both ergodic and several non-ergodic structures) and also for many other state of  
35 the art models. These new models provide state of the art results, while in the conference paper the  
36 results were not as good since they used an older, less sophisticated model.  
37
- 38 3. We have more experiments that demonstrate how learning complexity and scalability are handled  
39 by DHDPHMM versus HDPHMM and HMM.  
40
- 41 4. The phoneme recognition experiments in this paper are completely new.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



# A Left-to-Right HDP-HMM with HDPM Emissions

Amir Hossein Harati Nejad Torbati, Joseph Picone  
Department of Electrical and Computer Engineering  
Temple University  
Philadelphia, USA  
amir.harati@gmail.com, picone@temple.edu

Marc Sobel  
Department of Statistics  
Temple University  
Philadelphia, USA  
marc.sobel@temple.edu

**Abstract**— In this paper we introduce a new nonparametric Bayesian HMM based on the well-known HDP-HMM model. Unlike the original ergodic model, our model has a left-to-right structure. We introduce two approaches to adding non-emitting states that are used to model the beginning and end of finite duration sequences. Finally, we extend the HDP-HMM definition by introducing an HDP-HMM with HDP mixture emissions. We demonstrate that the new model outperforms the ergodic model for problems involving temporal structure by producing a 15% increase in likelihoods. Experiments on a phoneme classification task resulted in an 15.3% relative reduction in error.

**Keywords**—HDP-HMM; none-parametric Bayesian; Left-to-Right models; HMMs; Hierarchical Dirichlet Model

## I. INTRODUCTION

Hidden Markov models (HMMs) [1] are among the most powerful statistical modeling tools and have found a wide range of applications in many pattern recognition tasks such as speech recognition, machine vision, genomics and finance [2]. HMMs are parameterized both in their topology (e.g. number of states) and emission distributions (e.g. Gaussian mixtures). Model comparison methods are traditionally used to optimize the number of states and mixture components. However, these methods are computationally expensive and moreover there is no consensus on an optimum criterion for the selection [3].

An infinite HMM has been developed in the last few years [4][5][6] based on nonparametric Bayesian approaches. In this model, instead of defining a parametric prior over the transition distribution, a hierarchical Dirichlet process (HDP) prior is used. This model is known as an HDP-HMM model. HDP-HMM introduced in [5] and [6] is an ergodic model (a transition from an emitting state to all other states is allowed). However, in many pattern recognition applications involving temporal structure, such as speech processing, a left-to-right topology is preferred or sometimes required [7][8]. For example, in continuous speech recognition applications we model speech units (e.g. phonemes), which evolve in a sequential manner, using HMMs. Since we are dealing with an ordered sequence (e.g. a word is an ordered sequence of phonemes), a left-to-right model is preferred [7]. Moreover, the segmentation of speech data into these units is not known in advance, and therefore the training process must be able to connect these smaller models together into a larger HMM that

models the entire utterance. Obviously, this task can easily be achieved using left-to-right (LR) HMMs.

If the data has a finite length, the beginning and end of a sequence is typically modeled as two additional discrete events – non-emitting initial and final states [1][7]. In the original HDP-HMM formulation [5][6], this problem is not addressed. Also, the original HDP-HMM, as well as parametric HMMs, models each emission distribution by data points mapped to that state. For example, if we use a Gaussian mixture model (GMM) to model the emission distribution, for every state we compute a separate GMM and components can't be shared or re-used within a model. In this paper we propose a left-to-right HDP-HMM (LR HDP-HMM) with non-emitting initial and final states. In our model, emission distributions are modeled using GMMs with an infinite number of components. Sharing components is achieved by using an HDP prior instead of Dirichlet process (DP) priors as in [6].

The paper is organized as follows. In Section 2, we introduce Dirichlet processes and the HDP-HMM model. In Section 3, our proposed model is discussed. In Section 4, we present some experimental results on two datasets. We conclude the paper in Section 5 with a discussion of the limitations of the current model and future work.

## II. BACKGROUND

A Dirichlet process [9] is a discrete distribution that consists of countable infinite probability masses. A DP is denoted by  $DP(\alpha, H)$ , where  $\alpha$  is the concentration parameter and  $H$  is the base distribution. A DP can be represented by [10]:

$$G = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}, \quad \theta_k \sim H. \quad (1)$$

In this definition,  $\delta_{\theta_k}$  is the unit impulse function at  $\theta_k$ , and is referred to as an atom [5]. The weights  $\beta_k$  are sampled through a stick-breaking construction [5][10]:

$$\beta_k = v_k \prod_{l=1}^{k-1} (1 - v_l), \quad v_k | \alpha, G_0 \sim Beta(1, \alpha). \quad (2)$$

The sequence of  $\beta_k$  sampled by this process satisfies the constraint  $\sum_{k=1}^{\infty} \beta_k = 1$  with probability 1 and are denoted by  $\beta \sim GEM(\alpha)$  [5]. One of the main applications of a DP is to define a nonparametric prior distribution on the components of a mixture model. For example, a DP can be used to define a Gaussian mixture model (GMM) with an infinite number of mixture components [11]. This is a useful model in many areas of science. For example, in speech recognition, an acoustic unit (a word or a phoneme) can be modeled using a GMM [1].

A hierarchical Dirichlet process extends a DP to grouped data [5]. In this case there are several related groups and the goal is to model each group using a mixture model. These models can be linked using traditional parameter sharing approaches. For example, consider the problem of modeling acoustic units, such as phonemes, in continuous speech recognition using a mixture model in which parameters of different acoustic units can be shared. One approach is to use a DP to define a mixture model for each group and to use a global Dirichlet process,  $DP(\gamma, H)$ , as the common base distribution for all DPs [5]. An HDP is defined as:

$$\begin{aligned} G_0 &| \gamma, H \sim DP(\gamma, H) \\ G_j &| \alpha, G_0 \sim DP(\alpha, G_0), \end{aligned} \quad (3)$$

where  $H$  provides the prior for the parameters and  $G_0$  represents the average of the distribution of the parameters (e.g. means and covariances).

An alternative analogy, which is useful for gaining insight into the inference algorithms, is based on the concept of a Chinese restaurant franchise (CRF) [5]. In a CRF, a franchise consists of several restaurants with a common franchise-wide menu. Customers represent observed data, tables represent clusters and restaurants represent groups. The first customer entering restaurant  $j$  sits at one of the tables and orders an item from the menu. The next customer either sits at one of the occupied tables and eats the food served at that table or sits at a new table and orders new food from the menu. The probability of sitting at a table is proportional to the number of customers already seated at that table. However, if a customer starts a new table (with probability proportional to  $\alpha$ ), he or she orders food from the menu with a probability proportional to the number of tables serving that food in the franchise, or alternately orders a new food item with a probability proportional to  $\gamma$ .

An HDP-HMM [4][5][6] is an HMM with an unbounded number of states. In a typical ergodic HMM, the number of states is fixed so a matrix of dimension  $N$  states by  $N$  transitions per state is used to represent the transition probabilities. In an HDP-HMM, the transition matrix is replaced by an infinite, but discrete transition distribution, modeled by an HDP for each state. This lets each state have a different distribution for its transitions while the set of reachable states would be shared among all states. Fox et al. [6] extended the definition of HDP-HMM to HMMs with state persistence by introducing a sticky parameter  $\kappa$ . The definition for HDP-HMM is given by:

$$\begin{aligned} \beta &| \gamma \sim GEM(\gamma) \\ \pi_j &| \alpha, \beta \sim DP\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}\right) \\ \psi_j &| \sigma \sim GEM(\sigma) \\ \theta_{kj}^{**} &| H, \lambda \sim H(\lambda) \\ z_t &| z_{t-1}, \{\pi_j\}_{j=1}^{\infty} \sim \pi_{z_{t-1}} \\ s_t &| \{\psi_j\}_{j=1}^{\infty}, z_t \sim \psi_{z_t} \\ x_t &| \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t \sim F(\theta_{z_t, s_t}). \end{aligned} \quad (4)$$

The state, mixture component and observation are represented by  $z_t$ ,  $s_t$  and  $x_t$  respectively. The indices  $j$  and  $k$  are indices of the state and mixture components respectively. The base distribution that links all DPs together is represented by  $\beta$  and can be interpreted as the expected value of state transition distributions. The transition distribution for state  $j$  is a DP denoted by  $\pi_j$  with a concentration parameter  $\alpha$ . Another DP,  $\psi_j$ , with a concentration parameter  $\square$ , is used to model an infinite mixture model for each state ( $z_j$ ). The distribution  $H$  is the prior for the parameters  $\theta_{kj}$ . If we want the posterior distribution over the parameters to remain in the same family as the prior, then  $H$  should be chosen to be a conjugate prior to the observation likelihood. Since the likelihood has a multivariate normal distribution,  $H$  should have normal inverse Wishart (NIW) distribution.

### III. A LEFT-TO-RIGHT HDP-HMM WITH HDPM EMISSIONS

Hidden Markov models (HMMs) are a class of doubly stochastic processes in which discrete state sequences are modeled as a Markov chain [1]. The state of a Markov chain at time  $t$  is denoted by  $z_t$  and an observation is denoted by  $x_t \sim F(\theta_{z_t, s_t})$  where  $F$  is the emission distribution (e.g., a Gaussian mixture) and  $s_t$  is a mixture component index. In an HMM, there is a probability distribution to transit into state  $z_t$ . In an infinite HMM, this transition distribution should have infinite support and is modeled using HDP. For state  $j$  this transition distribution is denoted by  $\pi_j$ :

$$\pi_j | \alpha, \beta \sim DP\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_j}{\alpha + \kappa}\right). \quad (5)$$

From (5) we can see that the transition distribution has no topological restriction and therefore (4) defines an ergodic HMM. In this section we introduce a left-to-right HDP-HMM with initial and final non-emitting states. Moreover, we replace DP with HDP to model multimodal emission distributions that allow states to share mixture components.

#### A. Left-to-Right Transition Distributions

In order to obtain a left-to-right (LR) topology we need to force the base distribution of the Dirichlet distribution in (5) to only contain atoms to the right of the current state. This means  $\beta$  should be modified so that the probability of transiting to states left of the current state (i.e. states previously visited) becomes zero. For state  $j$  we define  $V_j = \{V_{ji}\}$  as:

$$V_{ji} = \begin{cases} 0, & i < j \\ 1, & i \geq j \end{cases} \quad (6)$$

where  $i$  is index for all states. Then we can modify  $\beta$  by multiplying it with  $V_j$ :

$$\beta' = \frac{\beta \cdot V_j}{\sum_i \beta_i V_{ji}}. \quad (7)$$

Therefore to obtain a left-to-right HDP-HMM, which we refer to as LR HDP-HMM, we simply replace  $\beta'$  with  $\beta$  in (5). The rest of the definition remains the same. Also notice that different topologies can be achieved by defining an appropriate  $V_j$ .

### B. Initial and Final Non-Emitting States

In many applications, such as continuous speech recognition, a LR HMM begins from and ends with non-emitting states. These states are required to model the beginning and end of finite duration sequences. Adding a non-emitting initial state is trivial: the probability of transition into the initial state is 1 and the probability distribution of a transition from this state is equal to  $\pi_{init}$  which is the initial probability distribution for an HDP-HMM without non-emitting states. However, adding a final non-emitting state is more complicated. In the following we will discuss two approaches to solving this problem.

#### 1) Maximum Likelihood Estimation

Consider state  $z_i$  depicted in Figure 1. The outgoing probabilities for any state can be classified into three categories: (1) a self-transition (P1), (2) a transition to all other states (P2), and (3) a transition to a final non-emitting state (P3). These probabilities must sum to 1:  $P1+P2+P3=1$ . Suppose that we obtained P2 from the inference algorithm. We will need to reestimate P1 and P3 from the data. This problem is, in fact, equivalent to the problem of tossing a coin until we obtain the first tails. Each head is equal to a self-transition and the first tails triggers a transition to the final state. This can be modeled using a geometric distribution [12]:

$$P(x = k) = (1 - \rho)^{k-1} \rho. \quad (8)$$

Equation (8) shows the probability of  $K-1$  heads before the first tail. In this equation  $1-\rho$  is the probability of heads (success). We also have:

$$\frac{P_1}{1-P_2} = 1-\rho, \quad \frac{P_3}{1-P_2} = \rho. \quad (9)$$

Suppose we have a total of  $N$  examples but for just  $M$  examples the state  $z_i$  is the last state of the model ( $S_M$ ). It can be shown [12] that the maximum likelihood estimation is obtained by:

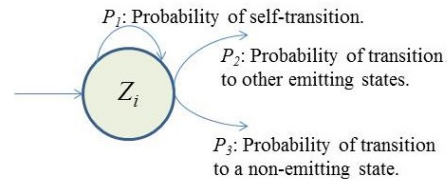


Figure 1- Outgoing probabilities for state  $z_i$

$$\hat{\rho} = \frac{M}{\sum_{i \in S_M} k_i} \quad (10)$$

where  $k_i$  are the number of self-transitions for state  $i$ . Notice that if  $z_i$  never happens to be the last state ( $M=0$ ),  $P3=0$ .

#### 2) Bayesian Estimation

Another approach to estimate  $\rho$  is to use a Bayesian framework. Since a beta distribution is the conjugate distribution for geometric distribution [13], we can use a beta distribution with hyperparameters  $(a,b)$  as the prior and obtain a posterior as [13][14]:

$$\rho \sim \text{Beta} \left( a + M, b + \sum_{i \in S_M} (k_i - 1) \right) \quad (11)$$

where  $M$  and  $S_M$  are same as in the previous section. Hyperparameters  $(a,b)$  can also be estimated using a Gibbs sampler if required [15].

### C. HDP Mixture Emission Distributions

In previous works [5][6], emission distributions for each state of an HDP-HMM were modeled using a Dirichlet process mixture (DPM) as shown in (4). While this model is reasonably flexible, each data point is strictly associated with a single state and hence statistical estimation of each parameter would be less reliable. This is a more serious problem for HDP-HMMs with a left-to-right topology since these models will discover more states. As a result the available data for estimating the emission distribution for each state would be more limited. The solution proposed here is to replace the DPM with an HDP mixture (HDPM) defined for the entire HMM. The final model without non-emitting states, which we refer to as LR HDP-HMM/HDPM, is defined by (12) and is displayed in Figure 2-(b). For comparison purposes, we display the original HDP-HMM in Figure 2-(a) [6].

$$\begin{aligned}
& \beta | \gamma \sim GEM(\gamma) \\
& \beta' = \frac{V_j \cdot \beta}{\sum_i V_{ji} \beta_i}, V_{ji} = \begin{cases} 0, & i < j \\ 1, & i \geq j \end{cases} \quad 1 \leq i < \infty \\
& \pi_j | \alpha, \beta' \sim DP(\alpha + \kappa, \frac{\alpha \beta' + \kappa \delta_j}{\alpha + \kappa}) \\
& \xi | \sigma \sim GEM(\sigma) \\
& \psi_j | \tau, \xi \sim DP(\tau, \xi) \\
& \theta_{kj}^{**} | H, \lambda \sim H(\lambda) \\
& z_t | z_{t-1}, \{\pi_j\}_{j=1}^{\infty} \sim \pi_{z_{t-1}} \\
& s_t | \{\psi_j\}_{j=1}^{\infty}, z_t \sim \psi_{z_t} \\
& x_t | \{\theta_{kj}^{**}\}_{k,j=1}^{\infty}, z_t \sim F(\theta_{z_t, s_t})
\end{aligned} \tag{12}$$

#### D. Modified Block Sampler

A block sampler for HDP-HMM with a multimodal emission distribution has been introduced by Fox et al. [6]. In this section we review the modifications of this algorithm needed for our new model. The interested reader should refer to [6][16] for additional details. The central idea is to jointly sample the state sequence  $z_{1:T}$  given the observations, model parameters and transition distribution  $\pi_j$ . A variant of forward-backward procedure [1] is utilized that allows us to exploit the Markovian structure of the HMM. However it requires approximation of the theoretically infinite distributions with a “degree  $L$  weak limit” approximation that truncates a DP into a Dirichlet distribution with  $L$  dimensions [17]:

$$GEM_L(\alpha) \triangleq Dir\left(\frac{\alpha}{L}, \dots, \frac{\alpha}{L}\right). \tag{13}$$

The sampling of the transition distribution is similar to [6]. The only difference is to replace  $\beta$  with  $\beta'$  given in (7). Using a similar approximation we can write the following prior distributions for the global weights  $\xi$  and state-specific weights  $\psi_j$  used in the HDPM emission distributions.

$$\xi | \sigma \sim Dir\left(\frac{\sigma}{L'}, \dots, \frac{\sigma}{L'}\right) \tag{14}$$

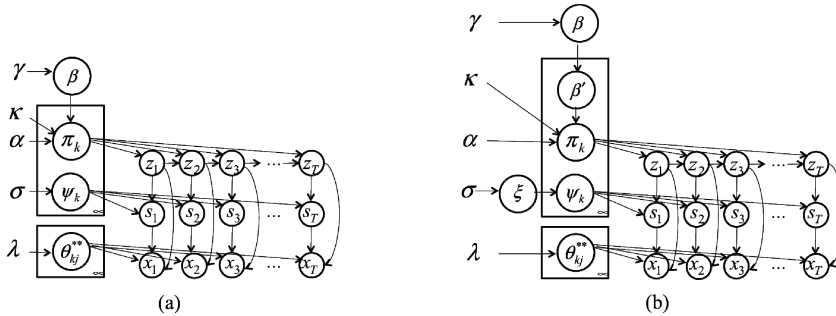


Figure 2- A comparison of models: (a) ergodic HDP-HMM [6] (b) proposed LR HDP-HMM/HDPM.

$$\psi_j | \xi, \tau \sim Dir(\tau \xi_1^j, \dots, \tau \xi_{L'}^j) \tag{15}$$

where  $L'$  is the order of approximation in this case. For the posterior distribution we can write:

$$\xi | M, \sigma \sim Dir\left(\frac{\sigma}{L'} + M_{\cdot 1}, \dots, \frac{\sigma}{L'} + M_{\cdot L'}\right) \tag{16}$$

$$\psi_j | \tau, \xi, Z_{1:T}, S_{1:T} \sim Dir(\tau \xi_1^j + n'_{j1}, \dots, \tau \xi_{L'}^j + n'_{jL'}) \tag{17}$$

where  $M_{jk}$  is the number of tables (clusters) in restaurant (state)  $j$  that serves dish (mixture component)  $k$ ;  $M_{\cdot k}$  is total number of tables in the franchise that serves dish  $k$ . The number of observations in state  $j$  that are assigned to component  $k$  is denoted by  $n'_{jk}$ . Estimating transition probabilities for the final non-emitting state can be done as a last step and after estimating the other parameters.

#### IV. EXPERIMENTS

**Synthetic data.** In the first experiment, we generate data from a left-to-right HMM without non-emitting states that consists of four states. The emission distribution for each state is a GMM with up to three components, each consisting of a two-dimensional normal distribution. Three synthetic data sequences totaling 1900 observations were generated for training. Three configurations have been studied: (1) an ergodic HDP-HMM, (2) a LR HDP-HMM with DPM emissions and (3) a LR HDP-HMM with HDPM emissions. An NIW prior is used for the mean and covariance. The truncation levels are set to 10 for both the number of states and the number of mixture components. Parameters of the NIW are set as follows: pseudocounts, the number of pseudo observations for the sample mean, is set to 0.1; the sample mean and covariance are set to the empirical mean and covariance; and degree of freedom, which is the precision on sample covariance, is set to 5.

Figure 3-(a) shows the average likelihoods for different models for held-out data by averaging five independent chains. Figure 3-(b) shows the structure of the models. The LR HDP-HMM/HDPM discovers the correct structure while the ergodic HDP-HMM finds a more simplified HMM. Moreover, we can see using HDP emissions improves the likelihood. While LR HDP-HMM/DPM can find the structure close to the correct one (not shown here), its likelihood is slightly less than that for the ergodic HDP-HMM. However, LR HDP-HMM/HDPM

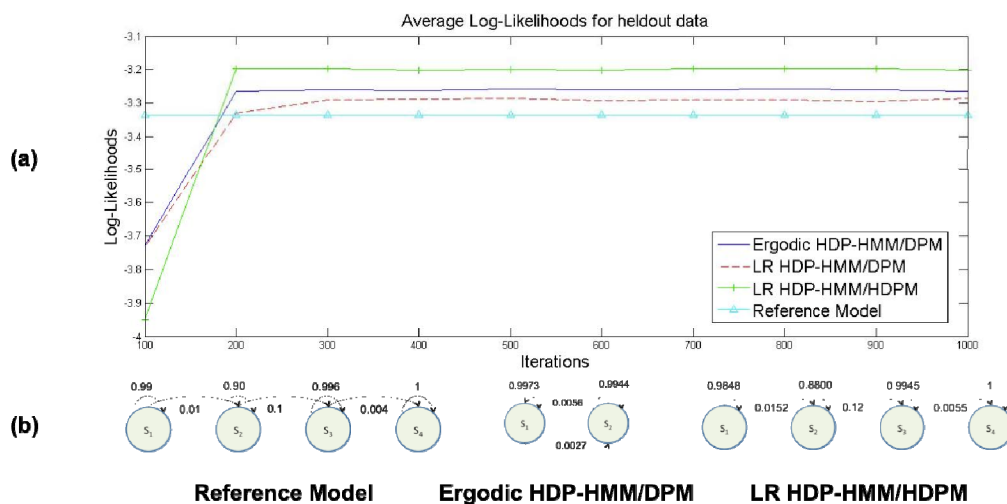


Figure 3- A comparison of (a) log-likelihoods of the proposed models to an ergodic model, and (b) the corresponding model structures.

produces a 15% improvement in likelihoods compared to the ergodic model. It is also interesting to note that the likelihoods of models discovered by all HDP-HMM algorithms are superior to the likelihood of the reference model itself.

**TIMIT Classification.** The TIMIT Corpus [18] is one of the most cited evaluation data sets used to compare new speech recognition algorithms. The data is segmented manually into phonemes and therefore is a natural choice to evaluate phoneme classification algorithms. TIMIT contains of 630 speakers from eight main dialects of American English [18]. The total numbers of utterances are 6300 where 3990 utterances are the standard training set and 150 utterances are core test set. We followed the standard practice of building models for 48 phonemes and then map them into 39 phonemes [19]. The first 12 Mel-Frequency Cepstral Coefficients (MFCCs) plus energy and their first and second derivatives features have been used to convert speech data into 39-dimensional feature streams. In this experiment, LR HDP-HMMs with Gaussian and DPM emissions have been used. We have used non-conjugate priors and placed a Gaussian prior on the mean and inverse-Wishart prior on the covariance matrix. Truncation levels are set to 10.

Table 1 compares the classification error of the left-to-right models and the parametric models. Since the maximum number of mixture components is set to 10, we have compared our systems to parametric HMMs with 10 components per

Table 1- A comparison of classification error rates

Model	Classification Error Rate
Parametric HMM [19] (10 mixtures)	27.8%
LR HDP-HMM with Gaussian emissions	26.7%
LR HDP-HMM with DPM emissions	24.1%

state. As this table shows, even left-to-right HDP-HMM with Gaussian emissions outperforms the parametric model.

Figure 4 shows the discovered structure for phonemes /aa/ and /sh/ using the proposed model. As the amount of data increases the system can learn a more complex model for the same phone. It is also important to note that the structure learned for each phone is different and reflects underlying differences between phonemes. Also note that the learned structure models multiple modalities by learning several parallel left-to-right paths. This is shown in Figure 4(c), where S1-S2, S1-S3 and S1-S4 depict three parallel models.

## V. CONCLUSION

In this paper we introduced a left-to-right HDP-HMM with HDPM emissions. We have shown that the new model can successfully learn the underlying structure when the data is generated using a generative left-to-right model. Moreover, it has been shown that the likelihood of the learned model is higher than the ergodic model. In this paper we have also introduced two approaches to adding non-emitting initial and final states to the left-to-right HDP-HMM model. Finally we presented the modifications needed in the block sampler to implement the inference algorithm for the new model. Through experimentation on TIMIT, we have shown that the proposed model outperforms parametric HMMs and can learn multimodal structure from the data.

One of the current problems of the HDP-HMM model (including left-to-right model) is that the inference algorithm is still computationally expensive. It is a serious problem when we are dealing with large datasets such as in speech or video processing applications. Therefore, our next task is to improve the inference algorithm specifically for left-to-right HDP-HMMs with HDPM emissions using its specific properties and structure. For example, due to the left-to-right constraints, the number of possible transitions in state 1 is L, in state 2 is L-1 and in state L is 1. We can exploit this fact to reduce the computational complexity.

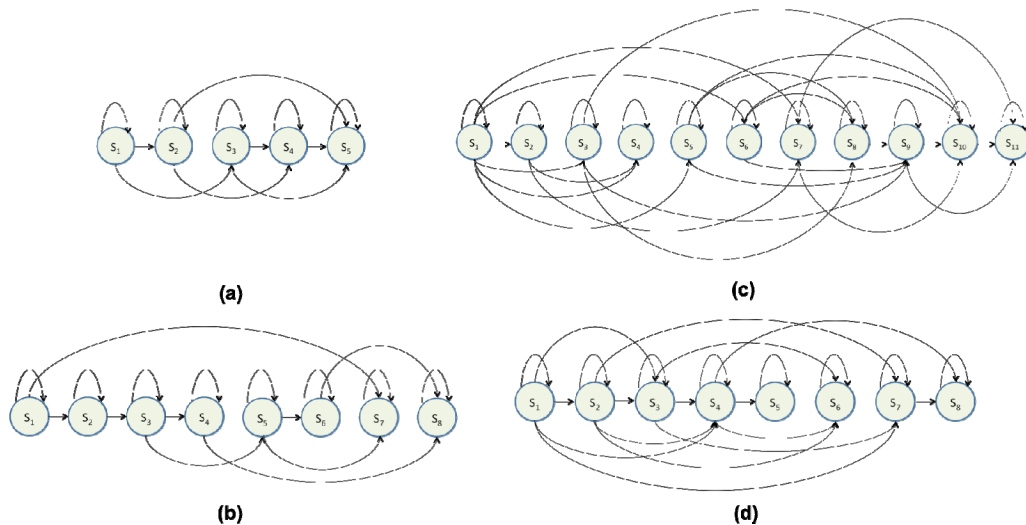


Figure 4- An automatically derived model structure for a left-to-right HDP-HMM model (without the first and last non-emitting states) for (a) /aa/ with 175 examples (b) /sh/ with 100 examples (c) /aa/ with 2,256 examples and (d) /sh/ with 1,317 examples. The data used in this illustration was extracted from the training portion of the TIMIT Corpus.

Another possible direction is to replace HDP emissions with more general hierarchical structures such as a Dependent Dirichlet Process [20] or an Analysis of Density (AnDe) model [21]. It has been shown that the AnDe model is the appropriate model for problems involves sharing statistical strength among multiple set of density estimators [5][21].

#### ACKNOWLEDGMENT

This research was supported in part by the National Science Foundation through Major Research Instrumentation Grant No. CNS-09-58854.

#### REFERENCES

- [1] L. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [2] P. Dymarski, *Hidden Markov Models, Theory and Applications*. InTech Open Access Publishers, 2011.
- [3] J. B. Kadane and N. A. Lazar, "Methods and Criteria for Model Selection," *Journal of the American Statistical Association*, vol. 99, no. 465, pp. 279–290, 2004.
- [4] M. Beal, Z. Ghahramani, and C. E. Rasmussen, "The Infinite Hidden Markov Model," in *Proceedings of Neural Information Processing Systems*, 2002, pp. 577–584.
- [5] Y. Teh, M. Jordan, M. Beal, and D. Blei, "Hierarchical Dirichlet Processes," *Journal of the American Statistical Association*, vol. 101, no. 47, pp. 1566–1581, 2006.
- [6] E. Fox, E. Sudderth, M. Jordan, and A. Willsky, "A Sticky HDP-HMM with Application to Speaker Diarization," *The Annals of Applied Statistics*, vol. 5, no. 2A, pp. 1020–1056, 2011.
- [7] B.-H. Juang and L. Rabiner, "Hidden Markov Models for Speech Recognition," *Technometrics*, vol. 33, no. 3, pp. 251–272, 1991.
- [8] G. A. Fink, "Configuration of Hidden Markov Models From Theory to Applications," in *Markov Models for Pattern Recognition*, Springer Berlin Heidelberg, 2008, pp. 127–136.
- [9] Y.-W. Teh, "Dirichlet process," in *Encyclopedia of Machine Learning*, Springer, 2010, pp. 280–287.
- [10] J. Sethuraman, "A constructive definition of Dirichlet priors," *Statistica Sinica*, vol. 4, no. 2, pp. 639–650, 1994.
- [11] C. E. Rasmussen, "The Infinite Gaussian Mixture Model," in *Proceedings of Advances in Neural Information Processing Systems*, 2000, pp. 554–560.
- [12] J. Pitman, *Probability*. Springer-Verlag, 1993.
- [13] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin, *Bayesian Data Analysis*, 2nd ed. Chapman & Hall, 2004.
- [14] P. Diaconis, K. Khare, and L. Saloff-Coste, "Gibbs Sampling, Conjugate Priors and Coupling," *Sankhya A*, vol. 72, no. 1, pp. 136–69, 2010.
- [15] F. A. Quintana and W. Tam, "Bayesian Estimation of Beta-binomial Models by Simulating Posterior Densities," *Journal of the Chilean Statistical Society*, vol. 13, no. 1–2, pp. 43–56, 1996.
- [16] E. Fox, E. Sudderth, M. Jordan, and A. Willsky, "Supplement to 'A Sticky HDP-HMM with Application to Speaker Diarization'," *The Annals of Applied Statistics*, vol. 5, no. 2A, pp. S1–S32, 2010.
- [17] H. Ishwaran and M. Zarepour, "Exact and approximate sum representations for the Dirichlet process," *Canadian Journal of Statistics*, vol. 30, no. 2, pp. 269–283, 2002.
- [18] J. Garofolo, L. Lamel, W. Fisher, J. Fiscus, D. Pallet, N. Dahlgren, and V. Zue, "TIMIT Acoustic-Phonetic Continuous Speech Corpus," *The Linguistic Data Consortium Catalog*, 1993.
- [19] A. Gunawardana, M. Mahajan, A. Acero, and J. C. Platt, "Hidden Conditional Random Fields for Phone Classification," in *Proceedings of INTERSPEECH*, 2005, pp. 1117–1120.
- [20] S. N. MacEachern, "Dependent Nonparametric Processes," in *ASA Proceedings of the Section on Bayesian Statistical Science*, 1999, pp. 50–55.
- [21] G. Tomlinson and M. Escobar, "Analysis of Densities," *Technical Report*, University of Toronto, Toronto, Canada, 1999.