

# Nonparametric Bayesian Models

## --Learning and Reasoning in Open Possible Worlds

Eric Xing

epxing@cs.cmu.edu

Machine Learning Dept./Language Technology Inst./Computer Science Dept.  
Carnegie Mellon University

## Outline

- Motivation and challenge
- Dirichlet Process and Infinite Mixture
  - Formulation
  - Approximate Inference algorithm
  - Example: population clustering
- Hierarchical Dirichlet Process and Multi-Task Clustering
  - Formulation
  - Transformed DP and HDP
  - Kernel stick-breaking process
  - Application: joint image segmentation
- Dynamic Dirichlet Process
  - Hidden Markov DP
  - Temporal DPM
  - Application: evolutionary clustering of documents
- Summary

# Clustering



# Image Segmentation



- How to segment images?
  - Manual segmentation (very expensive)
  - Algorithm segmentation
    - K-means
    - Statistical mixture models
    - Spectral clustering
- Problems with most existing algorithms
  - Ignore the spatial information
  - Perform the segmentation one image at a time
  - Need to specify the number of segments *a priori*

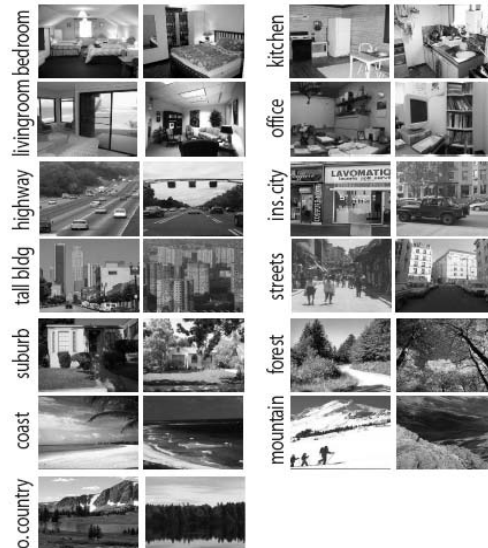
## Discover Object Categories

- Discover what objects are present in a collection of images in an unsupervised way

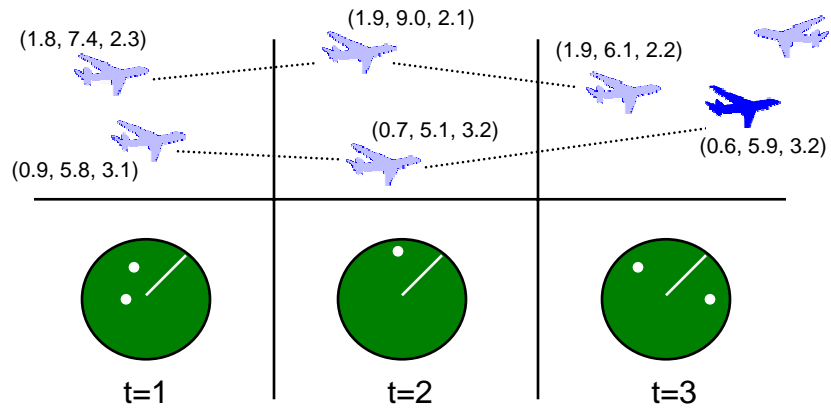


- Find those same objects in novel images
- Determine what local image features correspond to what objects; segmenting the image

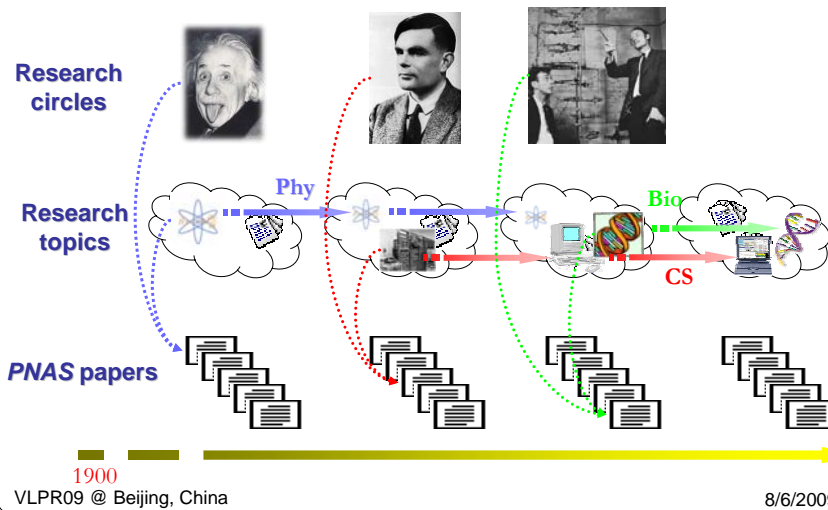
## Learn and Recognize Natural Scene Categories



# Object Recognition and Tracking

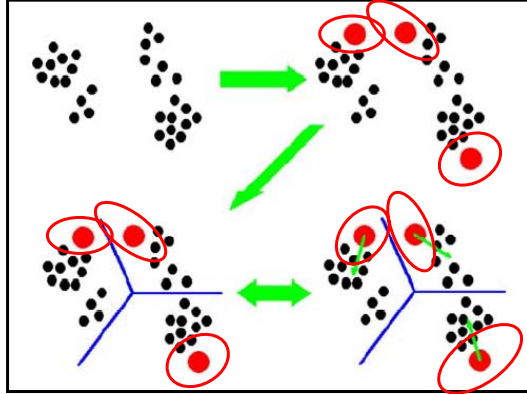


# The Evolution of Science



## A Classical Approach

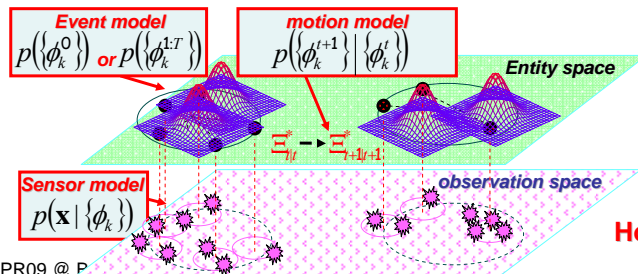
- Clustering as Mixture Modeling



- Then "model selection"

## Partially Observed, Open and Evolving Possible Worlds

- Unbounded # of objects/trajectories
- Changing attributes
- Birth/death, merge/split
- Relational ambiguity
- The parametric paradigm:



- Finite
- Structurally unambiguous

How to open it up?

## Model Selection vs. Posterior Inference

- Model selection
  - "intelligent" guess: ???
  - cross validation: data-hungry ☹
  - information theoretic:
    - AIC
    - TIC
    - MDL : }  $\arg \min KL(f(\cdot) | g(\cdot | \hat{\theta}_{ML}, K))$
  - Bayes factor: Parsimony, Ockam's Razor  
need to compute data likelihood
- Posterior inference:
  - we want to handle uncertainty of model complexity explicitly
  - $$p(M|D) \propto p(D|M)p(M)$$
  - $$M \equiv \{\theta, K\}$$
  - we favor a distribution that **does not** constrain  $M$  in a "closed" space!

## Two "Recent" Developments

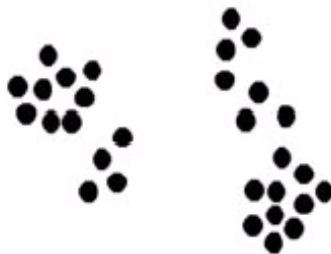
- First order probabilistic languages (FOPLs)
  - Examples: PRM, BLOG ...
  - Lift graphical models to "open" world (#rv, relation, index, lifespan ...)
  - Focus on complete, consistent, and operating rules to **instantiate** possible worlds, and formal language of expressing such rules
  - Operational way of defining distributions over possible worlds, via sampling methods
- Bayesian Nonparametrics
  - Examples: Dirichlet processes, stick-breaking processes ...
  - From finite, to infinite mixture, to more complex constructions (hierarchies, spatial/temporal sequences, ...)
  - Focus on the laws and behaviors of both the generative formalisms and resulting distributions
  - Often offer explicit expression of distributions, and expose the structure of the distributions --- motivate various approximate schemes

## Outline

- Motivation and challenge
- Dirichlet Process and Infinite Mixture
  - Formulation
  - Approximate Inference algorithm
  - Example: population clustering
- Hierarchical Dirichlet Process and Multi-Task Clustering
  - Formulation
  - Transformed DP and HDP
  - Kernel stick-breaking process
  - Application: joint image segmentation
- Dynamic Dirichlet Process
  - Hidden Markov DP
  - Temporal DPM
  - Application: evolutionary clustering of documents

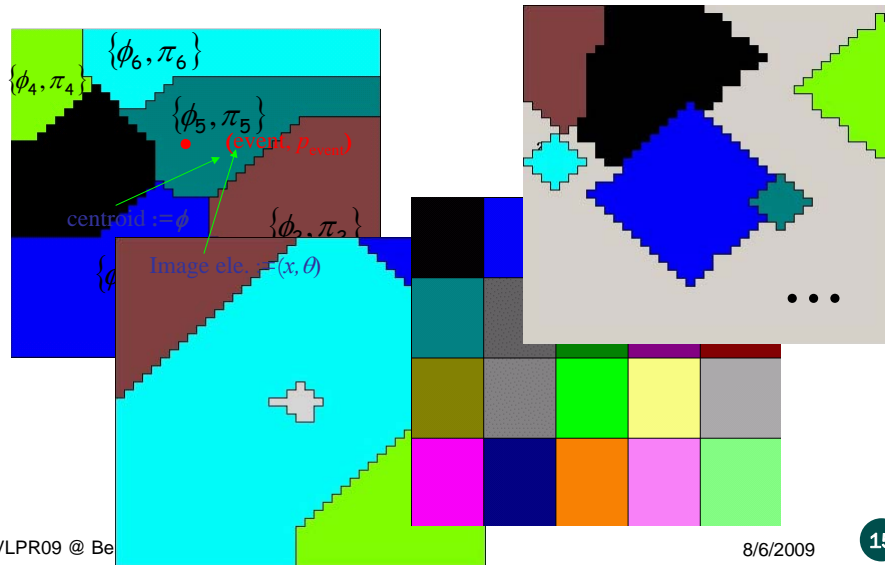
- Summary

## Clustering

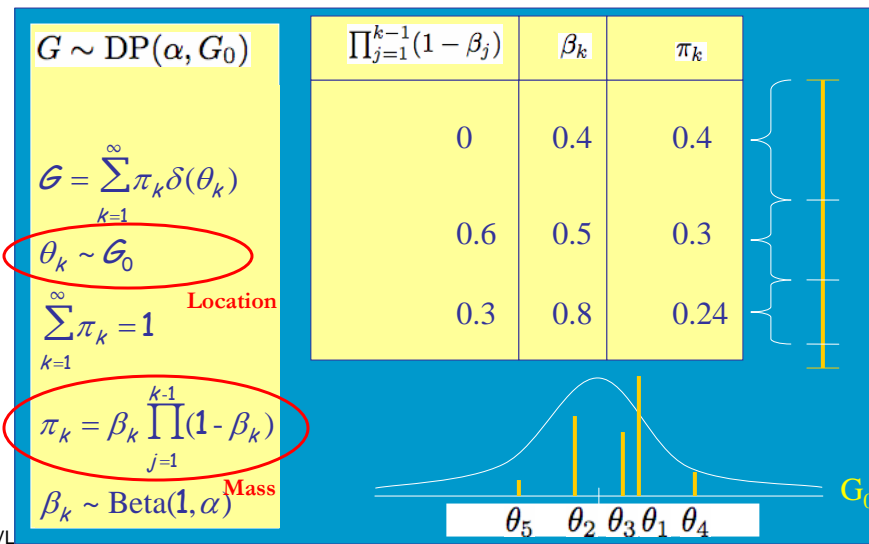


- How to label them ?
- How many clusters ???

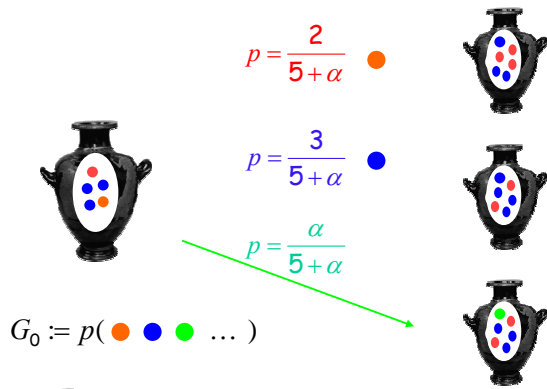
## Random Partition of Probability Space



## Stick-breaking Process



## DP – a Pólya urn Process

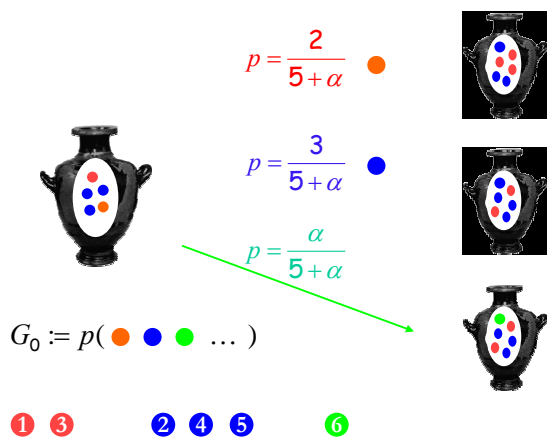


Joint:  $G(\text{urn}) \sim DP(\alpha G_0)$

Marginal:  $\phi_i | \phi_{-i}, \alpha, G_0 \sim \sum_{k=1}^K \frac{n_k}{i-1+\alpha} \delta_{\phi_k} + \frac{\alpha}{i-1+\alpha} G_0$

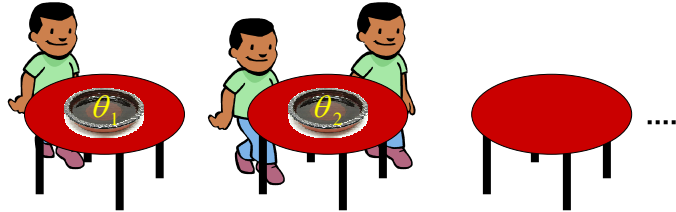
- Self-reinforcing property
- exchangeable partition of samples

## Clustering and DP Mixture



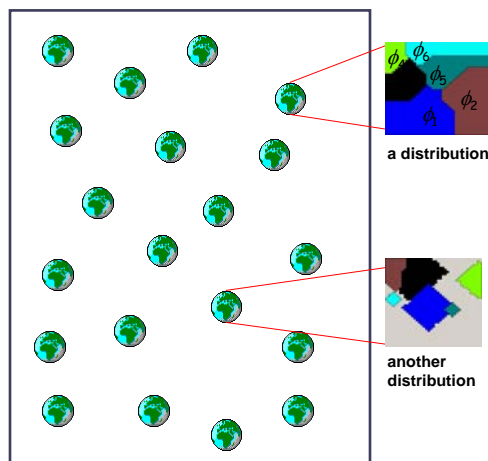
- We can associate mixture components with colors in the Pólya urn model and thereby define a *clustering* of the data

## Chinese Restaurant Process



$$P(c_i = k | \mathbf{c}_{-i}) = \begin{array}{ccc} 1 & 0 & 0 \\ \frac{1}{1+\alpha} & \frac{\alpha}{1+\alpha} & 0 \\ \frac{1}{2+\alpha} & \frac{1}{2+\alpha} & \frac{\alpha}{2+\alpha} \\ \frac{1}{3+\alpha} & \frac{2}{3+\alpha} & \frac{\alpha}{3+\alpha} \\ \frac{m_1}{i+\alpha-1} & \frac{m_2}{i+\alpha-1} & \dots \frac{\alpha}{i+\alpha-1} \end{array}$$

## Dirichlet Process



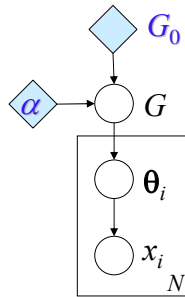
- A CDF,  $G$ , on possible worlds of random partitions follows a **Dirichlet Process** if for any measurable finite partition  $(\phi_1, \phi_2, \dots, \phi_m)$ :

$$(G(\phi_1), G(\phi_2), \dots, G(\phi_m)) \sim \text{Dirichlet}(\alpha G_0(\phi_1), \dots, \alpha G_0(\phi_m))$$

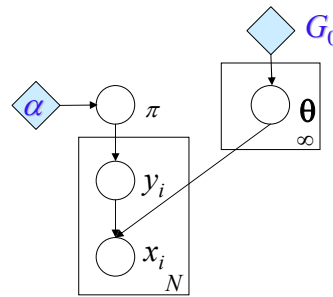
where  $G_0$  is the **base measure** and  $\alpha$  is the **scale parameter**

Thus a Dirichlet Process  $G$  defines a distribution of distribution

# Graphical Model Representations of DP



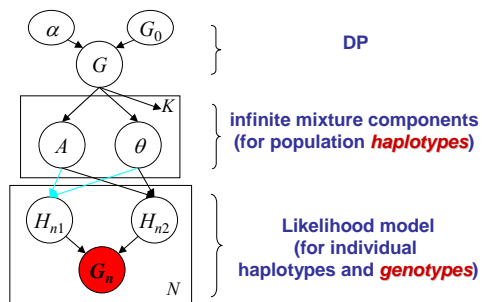
The Pólya urn construction



The Stick-breaking construction

# Example: DP-haplotype [Xing et al, 2004]

- Clustering human populations



- Inference: Markov Chain Monte Carlo (MCMC)
  - Gibbs sampling
  - Metropolis Hasting

## Inheritance and Observation Models

- Single-locus mutation model

$$A_{C_{i_e}} \rightarrow H_{i_e}$$

$$P_H(h_i | a_i, \theta) = \begin{cases} \theta & \text{for } h_i = a_i \\ \frac{1-\theta}{|B|-1} & \text{for } h_i \neq a_i \end{cases}$$

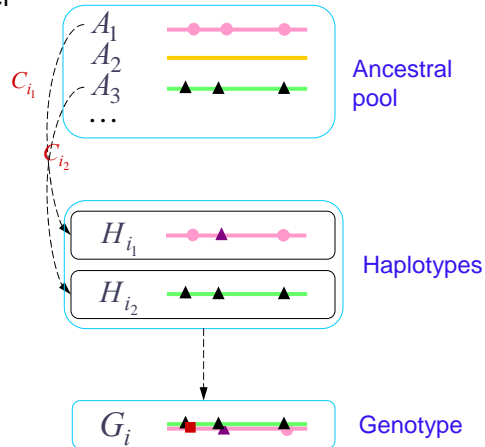
$\rightarrow h_i = a_i$  with prob.  $\theta$

- Noisy observation model

$$H_{i_1}, H_{i_2} \rightarrow G_i$$

$$P_G(g | h_1, h_2):$$

$$g_i = h_{1,i} \oplus h_{2,i} \text{ with prob. } \lambda$$



## MCMC for Haplotype Inference

- Gibbs sampling for exploring the posterior distribution under the proposed model
  - Integrate out the parameters such as  $\theta$  or  $\lambda$ , and sample  $c_{i_e}$ ,  $a_k$  and  $h_{i_e}$

$$p(c_{i_e} = k | \mathbf{c}_{[-i_e]}, \mathbf{h}, \mathbf{a}) \propto p(c_{i_e} = k | \mathbf{c}_{[-i_e]}) p(h_{i_e} | a_k, \mathbf{h}_{[-i_e]}, \mathbf{c})$$

Posterior                      Prior                      x                      Likelihood

⋮

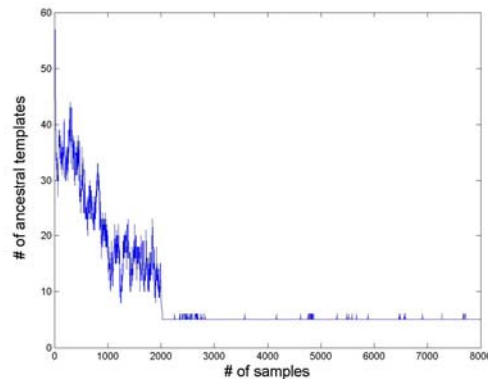
Pólya urn

- Gibbs sampling algorithm: draw samples of each random variable to be sampled given values of all the remaining variables

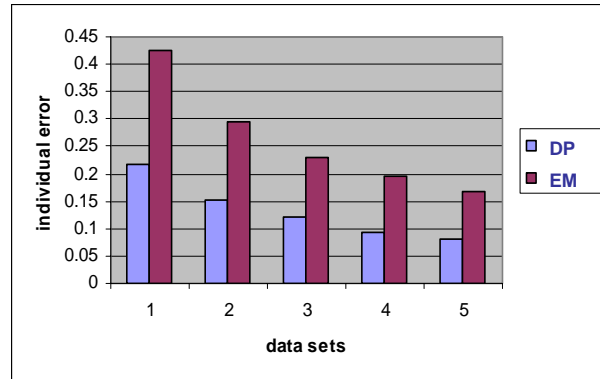
## MCMC for Haplotype Inference

1. Sample  $c_{ie}^{(j)}$ , from  $p(c_{ie}^{(j)} = k | \mathbf{c}^{[-j, ie]}, \mathbf{h}, \mathbf{a})$   
 $\propto p(c_{ie}^{(j)} = k | \mathbf{c}^{[-j, ie]}, \mathbf{m}, \mathbf{n}) p(h_{ie}^{(j)} | a_k, \mathbf{c}, \mathbf{h}^{[-j, ie]})$   
 $\propto (m_{jk}^{[-j, ie]} + \tau \beta_k) p(h_{ie}^{(j)} | a_k, l_k^{[-j, ie]}), \text{ for } k = 1, \dots, K + 1$
2. Sample  $a_k$  from  $p(a_{k,t} | \mathbf{c}, \mathbf{h}) \propto \prod_{j, ie | c_{ie}^{(j)} = k} p(h_{ie,t}^{(j)} | a_{k,t}, l_{k,t}^{(j)})$   
 $= \frac{\Gamma(\alpha_h + l_{k,t}) \Gamma(\beta_h + l'_{k,t})}{\Gamma(\alpha_h + \beta_h + m_k)} R(\alpha_h, \beta_h)$
3. Sample  $h_{ie}^{(j)}$  from  $p(h_{ie,t}^{(j)} | \mathbf{h}_{[-ie,t]}^{(j)}, \mathbf{c}, \mathbf{a}, \mathbf{g})$ 
  - For DP scale parameter  $\alpha$ : a vague inverse Gamma prior

## Convergence of Ancestral Inference



## DP vs. Finite Mixture via EM



## Variational Inference [Blei & Jordan 2005, Kurihara et al 2007]

- Gibbs sampling solution is not efficient enough to scale up to the large scale problems.
- Truncated stick-breaking approximation can be formulated in the space of explicit, non-exchangeable cluster labels.
- Variational inference can now be applied to such a finite-dimensional distribution
- Variational Inference:
  - For a complicated  $P(X_1, X_2, \dots, X_n)$ , approximate it with  $Q(X)$ :

$$Q(\mathbf{X}) = \prod_i Q(\mathbf{X}_{C_i})$$

$$\{Q^*(\mathbf{X}_{C_i})\} = \arg \min KL(Q(\mathbf{X})|P(\mathbf{X}))$$

# Approximations to DP

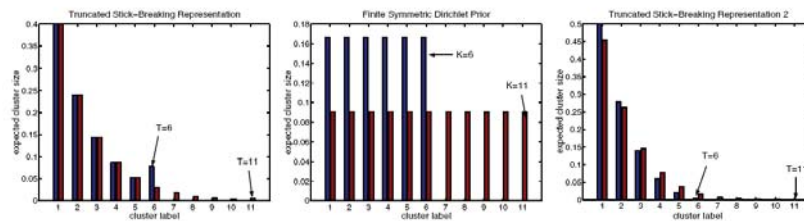
- Truncated stick-breaking representation
- Finite symmetric Dirichlet approximation

$$\begin{aligned}
 v_i &\sim \mathcal{B}(v_i; 1, \alpha) & i = 1, \dots, T-1 & & \pi &\sim \mathcal{D}(\pi; \frac{\alpha}{K}, \dots, \frac{\alpha}{K}) \\
 v_T &= 1 & & & & \\
 \pi_i &= v_i \prod_{j < i} (1 - v_j) & i = 1, \dots, T & & & \\
 \pi_i &= 0 & i > T & & & 
 \end{aligned}$$

- The joint distribution can be expressed as:
- The joint distribution can be expressed as:

$$\begin{aligned}
 P(X, z, v, \eta) &= \left[ \prod_{n=1}^N p(x_n | \eta_{z_n}) p(z_n | \pi(v)) \right] \left[ \prod_{i=1}^T p(\eta_i) \mathcal{B}(v_i; 1, \alpha) \right] \\
 P(X, z, \pi, \eta) &= \left[ \prod_{n=1}^N p(x_n | \eta_{z_n}) p(z_n | \pi) \right] \left[ \prod_{i=1}^K p(\eta_i) \right] \mathcal{D}(\pi; \frac{\alpha}{K}, \dots, \frac{\alpha}{K})
 \end{aligned}$$

# TDP vs. TSB



- TDP is size biased
- cluster labels is NOT interchangeable under TDP but is interchangeable under TSB

## Marginalization

$$P(\mathbf{X}, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\eta}) = \left[ \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\eta}_{z_n}) p(z_n | \boldsymbol{\pi}) \right] \left[ \prod_{i=1}^K p(\eta_i) \right] \mathcal{D}(\boldsymbol{\pi}; \frac{\alpha}{K}, \dots, \frac{\alpha}{K})$$

- In variational Bayesian approximation, we assume a factorized form for the posterior distribution.
- However it is not a good assumption since changes in  $\boldsymbol{\pi}$  will have a considerable impact on  $\mathbf{z}$ .

If we can integrate out  $\boldsymbol{\pi}$ , the joint distribution is given by

$$P(\mathbf{X}, \mathbf{z}, \boldsymbol{\eta}) = \left[ \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\eta}_{z_n}) \right] p(\mathbf{z}) \left[ \prod_{i=1}^{\infty} p(\eta_i) \right]$$

For the TSB representation:

$$p_{\text{TSB}}(\mathbf{z}) = \prod_{i < T} \frac{\Gamma(1 + N_i) \Gamma(\alpha + N_{>i})}{\Gamma(1 + \alpha + N_{\geq i})} \alpha$$

For the FSD representation:

$$p_{\text{FSD}}(\mathbf{z}) = \frac{\Gamma(\alpha) \prod_{k=1}^K \Gamma(N_k + \frac{\alpha}{K})}{\Gamma(N + \alpha) \Gamma(\frac{\alpha}{K})^K}$$

## VB inference

- We can then apply the VB inference on the four approximations

$$\{Q^*(\mathbf{X}_{C_i})\} = \arg \min KL(Q(\mathbf{X}) | P(\mathbf{X}))$$

The approximated posterior distribution for TSB and FSD are

$$Q_{\text{TSB}}(\mathbf{z}, \boldsymbol{\eta}, \mathbf{v}) = \left[ \prod_n^N q(z_n) \right] \left[ \prod_{i=1}^T q(\eta_i) q(v_i) \right] \quad Q_{\text{FSD}}(\mathbf{z}, \boldsymbol{\eta}, \boldsymbol{\pi}) = \left[ \prod_n^N q(z_n) \right] \left[ \prod_{k=1}^K q(\eta_k) \right] q(\boldsymbol{\pi})$$

Depending on marginalization or not,  $\mathbf{v}$  and  $\boldsymbol{\pi}$  may be integrated out.

## Experimental results

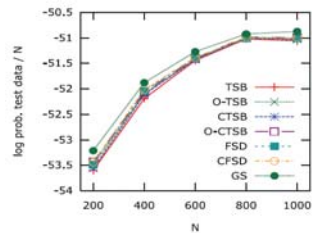


Figure 2: Average log probability per data-point for test data as a function of  $N$ .

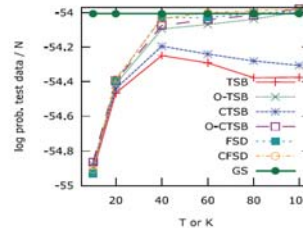
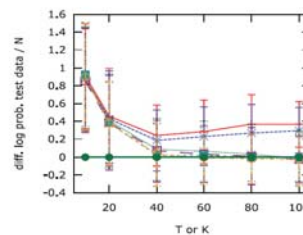
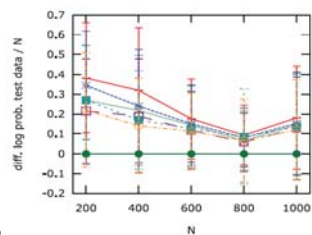


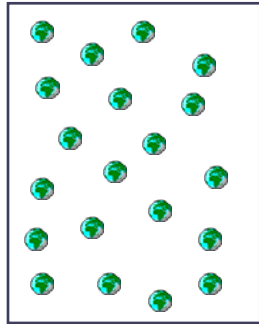
Figure 4: Average log probability per data-point for test data as a function of  $T$  (for TSB methods) or  $K$  (for FSD methods).



## Outline

- Motivation and challenge
- Dirichlet Process and Infinite Mixture
  - Formulation
  - Approximate Inference algorithm
  - Example: population clustering
- Hierarchical Dirichlet Process and Multi-Task Clustering
  - Formulation
  - Transformed DP and HDP
  - Kernel stick-breaking process
  - Application: joint image segmentation
- Dynamic Dirichlet Process
  - Hidden Markov DP
  - Temporal DPM
  - Application: evolutionary clustering of documents
- Summary

## Solving Multiple Clustering Problems

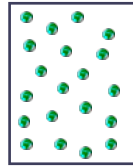


$$G(\Psi) \sim DP(\alpha G_0)$$

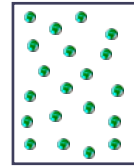
1 3

2 4 5

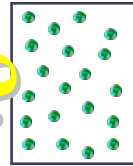
6



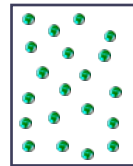
$$G(\Psi) \sim DP(\alpha G_0)$$



$$G(\Psi) \sim DP(\alpha G_0)$$



$$G(\Psi) \sim DP(\alpha G_0)$$



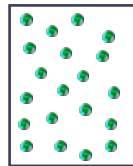
$$G(\Psi) \sim DP(\alpha G_0)$$

8/6/2009

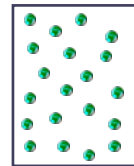
35

## Solving Multiple Clustering Problems

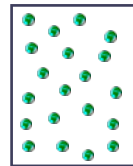
- Solve separately
  - Fail to capture correlation
  - Fail to cross-reinforce shared information (i.e., topic specific lexicon)
  - Data fragmentation



$$G(\Psi) \sim DP(\alpha G_0)$$



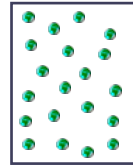
$$G(\Psi) \sim DP(\alpha G_0)$$



$$G(\Psi) \sim DP(\alpha G_0)$$



- Solve together
  - Then what is the difference between all these journals?



$$G(\Psi) \sim DP(\alpha G_0)$$

8/6/2009

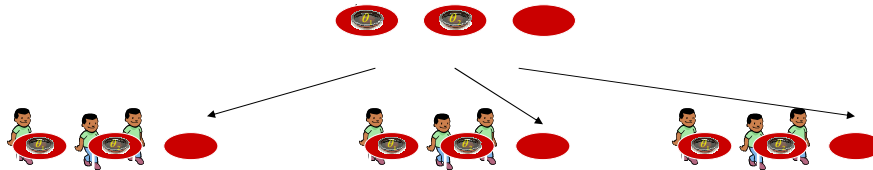
36

# Hierarchical Dirichlet Process

[Teh et al., 2005, Xing et al. 2005]

- Two level Pólya urn scheme

- At the  $i$ -th step in  $j$ -th "group",



- Choose  $\theta_k$  with prob.  $\frac{m_{jk}}{\sum_k m_{jk} + \alpha_0}$
- Go to the upper level DP with prob.  $\frac{\alpha_0}{\sum_k m_{jk} + \alpha_0}$

Oracle

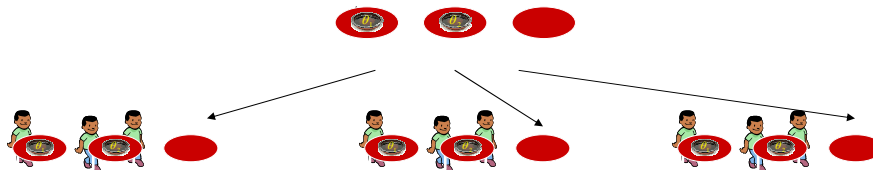
- Choose  $\theta_k$  with prob.  $\frac{n_k}{\sum n_k + \gamma}$
- Draw a new sample with prob.  $\frac{\gamma}{\sum n_k + \gamma}$

# Hierarchical Dirichlet Process

[Teh et al., 2005, Xing et al. 2005]

- Two level Pólya urn scheme

- At the  $i$ -th step in  $j$ -th "group",



- Draw from stock urn define Dirichlet Process  $DP(\gamma, H)$

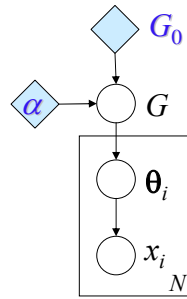
$$\theta_i | \theta_{-i} \sim \sum_{k=1}^K \frac{n_k}{i + \gamma} \delta_{\theta_k}(\theta_i) + \frac{\gamma}{i + \gamma} H(\theta_i)$$

- Conditioning on  $DP(\gamma, H)$ , the  $m$ th draw from the  $m$ th bottom-level urn also form a Dirichlet measure

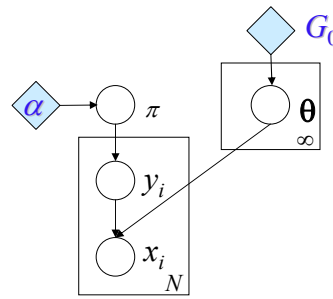
$$\theta_{m_j} | \theta_{-m_j} \sim \sum_{k=1}^K \frac{m_{j,k} + \alpha}{m_j + \alpha} \frac{n_k}{n + \gamma} \delta_{\theta_k}(\theta_{m_j}) + \frac{\alpha}{m_j + \alpha} \frac{\gamma}{i + \gamma} H(\theta_{m_j})$$

$$= \sum_{k=1}^K p_k^j \delta_{\theta_k}(\theta_{m_j}) + p_{k+1}^j H(\theta_{m_j})$$

# Recall: Graphical Model Representations of DP

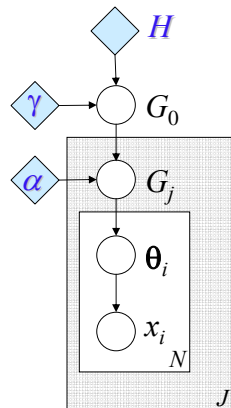


The Pólya urn construction



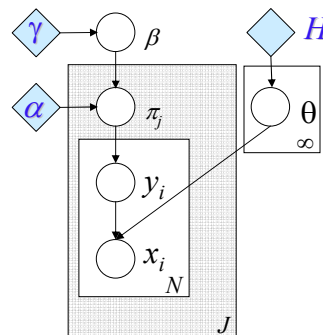
The Stick-breaking construction

# Hierarchical DP Mixture



Stick( $\alpha, \beta$ ):

$$\pi_{jk} \sim \text{Beta}(\alpha\beta_k, \alpha(1 - \sum_{l=1}^k \beta_l)), \quad \pi_{jk} = \pi_{j\cdot} \prod_{l=1}^{k-1} (1 - \pi'_{jl}).$$

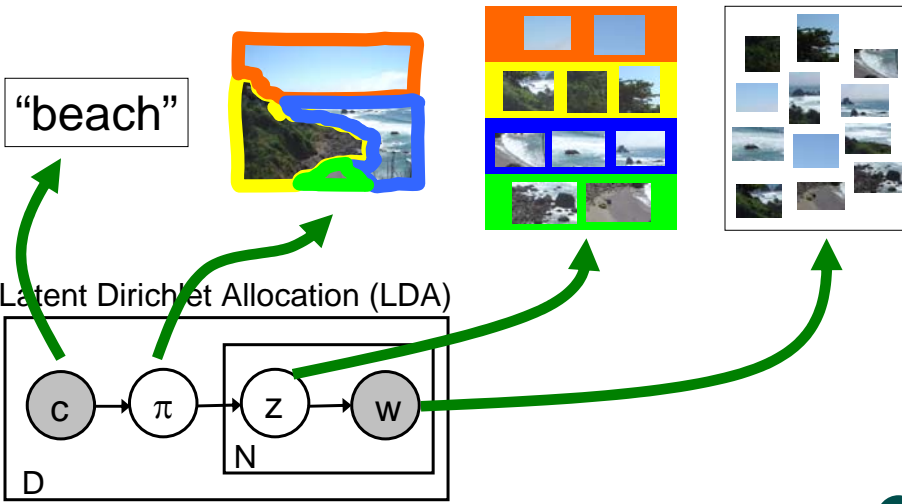


$\theta_k \sim H$

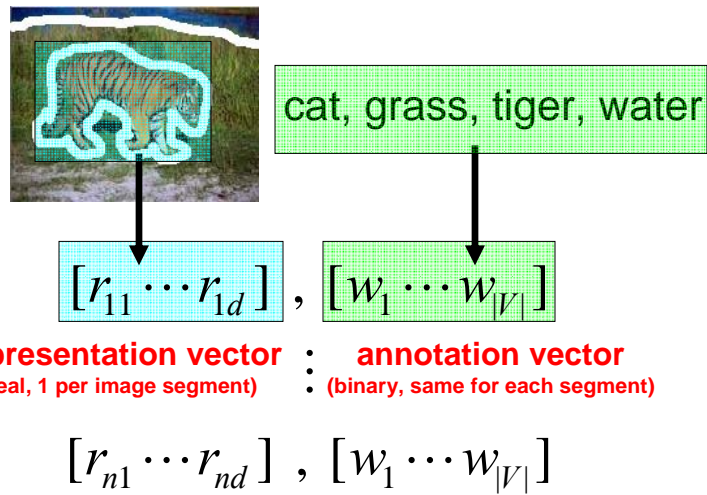
$$\beta = \text{Stick}(\gamma), \quad G_0 = \sum_{k=1}^{\infty} \beta_k \delta(\theta_k)$$

$$\pi_j = \text{Stick}(\alpha, \beta), \quad G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta(\theta_k)$$

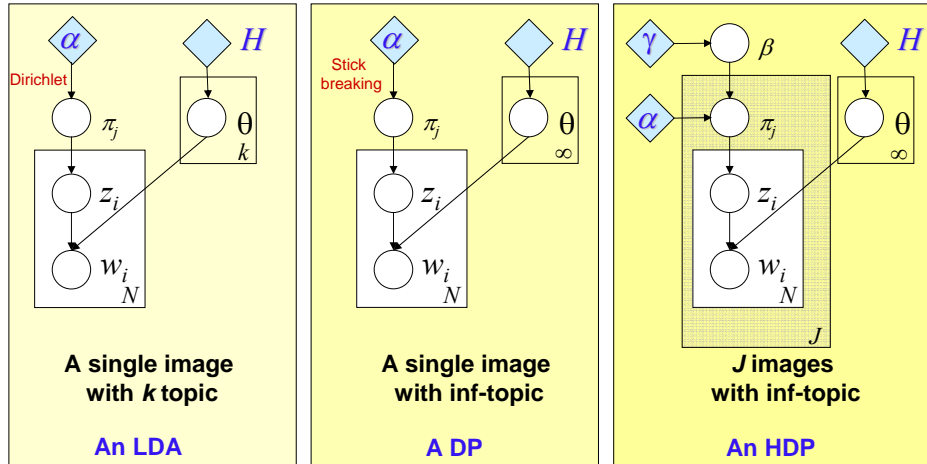
## Topic Models for Images



## Image Representation

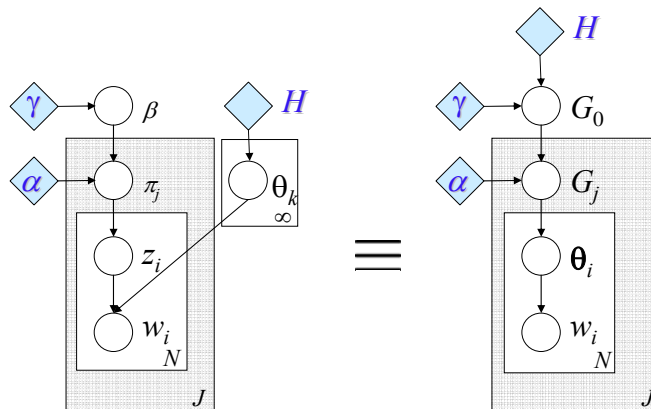


# Infinite Topic Model for Image



# Problem with HDP

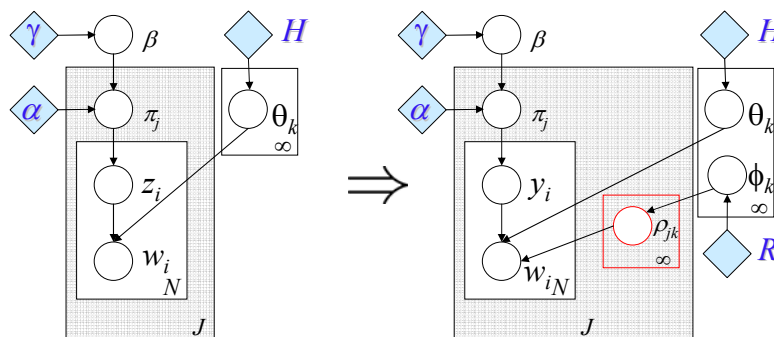
- Every group (i.e., image) has exactly the same set of visual-vocabulary topics, albeit with different frequency



## Transformed Dirichlet Process

[Sudderth et al, 2005]

- An extension of HDP in which global mixture components undergo a set of random transformations before being reused in each group



## Synthetic Data Results

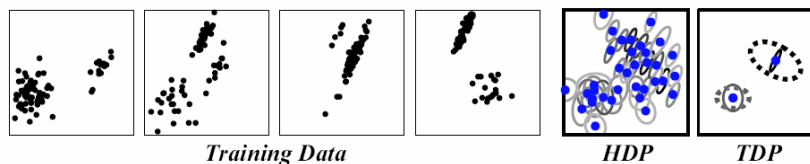
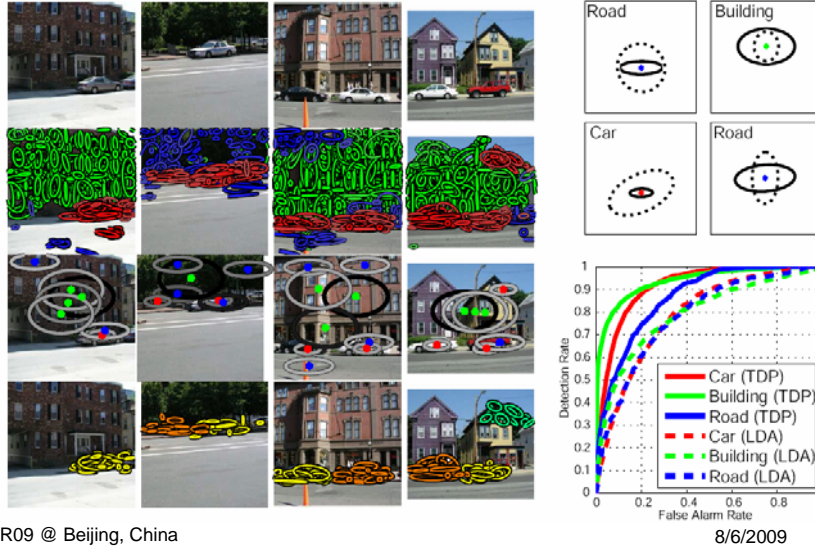


Figure 3: Comparison of hierarchical models learned via Gibbs sampling from synthetic 2D data. *Left*: Four of 50 “images” used for training. *Center*: Global distribution  $G_0(\theta)$  for the HDP, where ellipses are covariance estimates and intensity is proportional to prior probability. *Right*: Global TDP distribution  $G_0(\theta, \rho)$  over both clusters  $\theta$  (solid) and translations  $\rho$  of those clusters (dashed).

- HDP uses a large set of global clusters to discretize the transformations underlying the data, and may have poor generalization for modeling visual scenes.

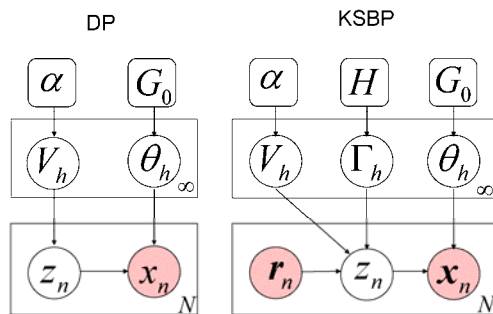
## Analyzing Street Scenes



## Kernel stick-breaking process

[Dunson and Park, 2006]

- For image analysis, we want to impose the belief that spatially proximate patches are more probable to be associated with the same cluster.
- We augmented the stick-breaking representation of DP to employ a kernel function to quantify some additional prior.

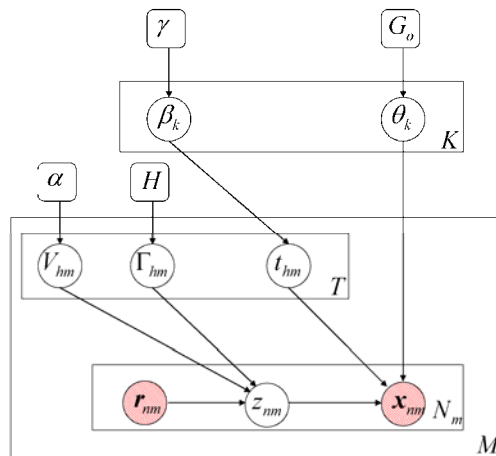


## KSBP for image analysis

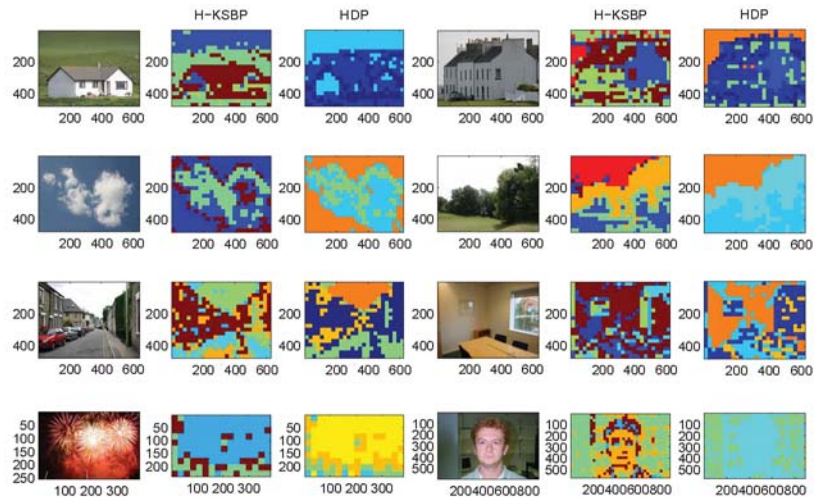
- $N$  • Consider a image composed of  $N$  patches, the features vectors  $\{\mathbf{x}_n\}_{n=1,N}$  and the associated locations  $\{\mathbf{r}_n\}_{n=1,N}$  can be modeled as follows:

$$\begin{aligned} \mathbf{x}_n &\sim f(\phi_n) \text{ iid} \\ \phi_n &\sim G_r \text{ iid} \\ G_r &= \sum_{h=1}^{\infty} \pi_h(\mathbf{r}; V_h, \Gamma_h, \psi) \delta_{\theta_h} \\ \pi_h(\mathbf{r}; V_h, \Gamma_h, \psi) &= V_h K(\mathbf{r}, \Gamma_h, \psi) \prod_{l=1}^{h-1} [1 - V_l K(\mathbf{r}, \Gamma_l, \psi)] \\ V_h &\sim \text{Beta}(a, b) \text{ iid} \\ \Gamma_h &\sim H \text{ iid} \\ \theta_h &\sim G_o \text{ iid} \end{aligned}$$

## Multi-task Image Segmentation



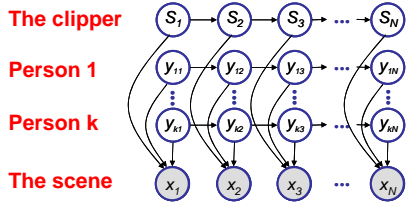
## Segmentation results [An et al, 2008]



## Outline

- Motivation and challenge
- Dirichlet Process and Infinite Mixture
  - Formulation
  - Approximate Inference algorithm
  - Example: population clustering
- Hierarchical Dirichlet Process and Multi-Task Clustering
  - Formulation
  - Transformed DP and HDP
  - Kernel stick-breaking process
  - Application: joint image segmentation
- Dynamic Dirichlet Process
  - Hidden Markov DP
  - Temporal DPM
  - Application: evolutionary clustering of documents
- Summary

# Object Recognition and Tracking



- Each chain corresponds to the trajectory of a specific object



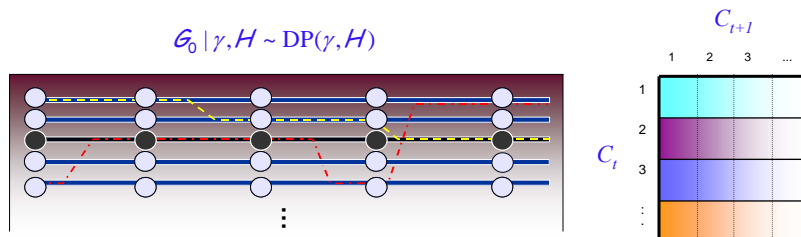
Multi-View Face Tracking with Factorial and Switching HMM

Peng Wang, Qiang Ji  
Department of Electrical, Computer and System Engineering  
Rensselaer Polytechnic Institute  
Troy, NY 12180

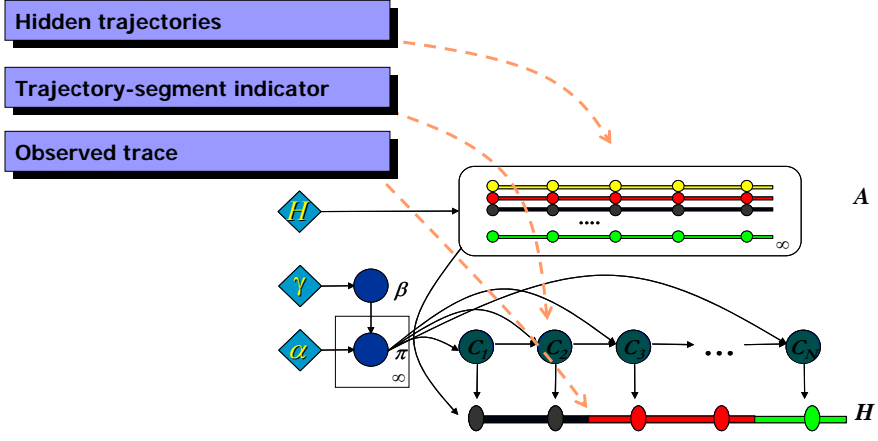
# Hidden Markov Dirichlet Process

(Xing and Sohn. Bayesian Analysis, 2007, Sohn and Xing, ISMB 2007)

- Hidden Markov Dirichlet process mixtures
  - Extension of HMM model to infinite ancestral space
    - Infinite dimensional transition matrix
    - Each row of the transition matrix is modeled with a DP:  $G_i | \alpha, G_0 \sim DP(\alpha, G_0)$

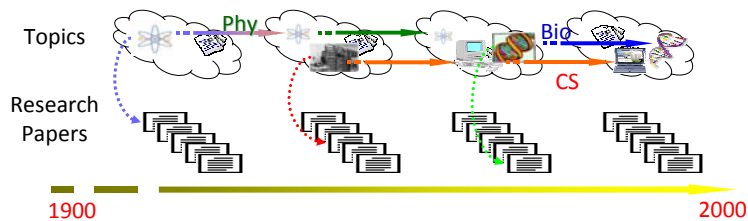


# HMDP as a Graphical Model



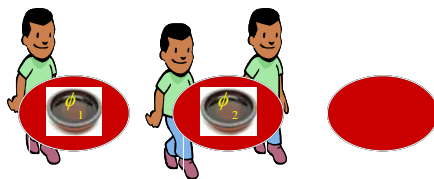
# Evolutionary Clustering

- Adapts the number of mixture components over time
  - Mixture components can die out
  - New mixture components are born at any time
  - Retained mixture components parameters evolve according to a Markovian dynamics



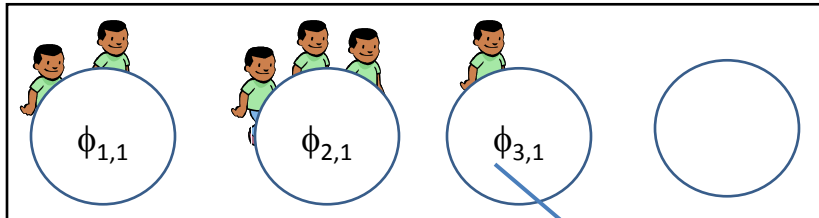
## The Chinese Restaurant Process

- Customers correspond to **data points**
- Tables correspond to **clusters/mixture components**
- Dishes correspond to **parameter** of the mixtures



## Temporal DPM [Ahmed and Xing 2008]

- **The Recurrent Chinese Restaurant Process**
  - The restaurant operates in **epochs**
  - The restaurant is **closed** at the end of each epoch
  - The **state** of the restaurant at time epoch  $t$  **depends** on that at time epoch  $t-1$ 
    - Can be extended to **higher-order** dependencies.

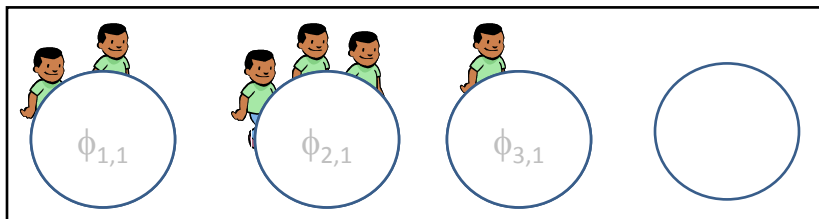


T=1

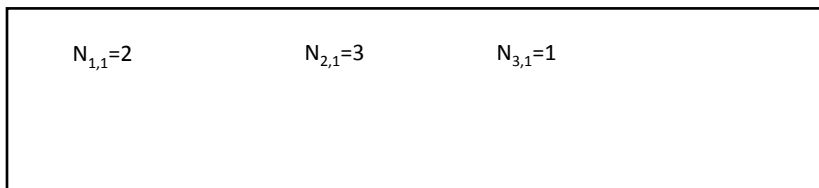
Dish eaten at table 3 at time epoch 1  
OR the parameters of cluster 3 at time epoch 1

#### Generative Process

- Customers at time  $T=1$  are seated as before:
  - Choose table  $j \propto N_{j,1}$  and Sample  $x_i \sim f(\phi_{j,1})$
  - Choose a new table  $K+1 \propto \alpha$
  - Sample  $\phi_{K+1,1} \sim G_0$  and Sample  $x_i \sim f(\phi_{K+1,1})$

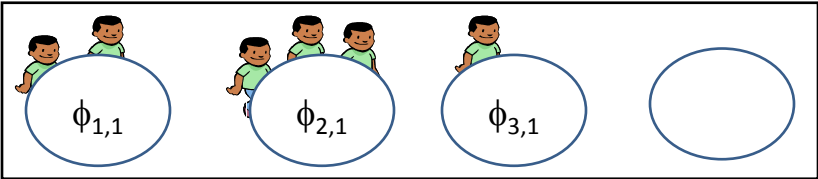


T=1

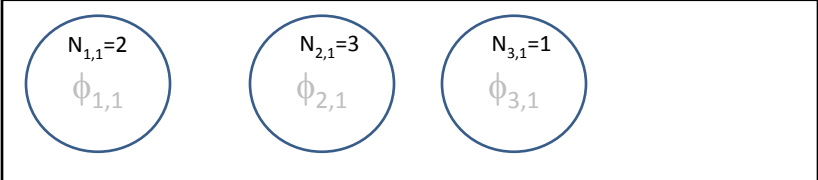


T=2


**SAILING LAB**  
Laboratory for Multiscale Analysis and Optimal Management



T=1

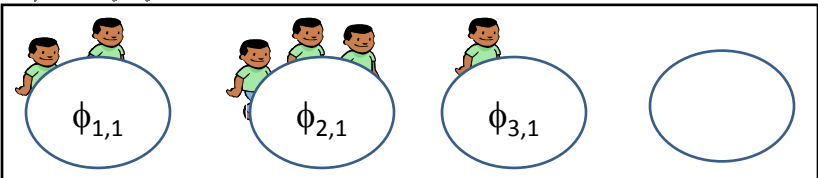


T=2

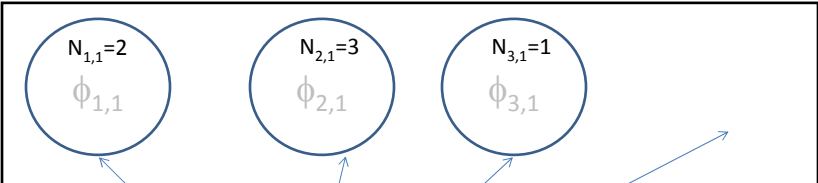


VLPR09 @ Beijing, China
8/6/2009
61


**SAILING LAB**  
Laboratory for Multiscale Analysis and Optimal Management



T=1



T=2



$\frac{2}{6+\alpha}$

$\frac{3}{6+\alpha}$

$\frac{1}{6+\alpha}$

$\frac{\alpha}{6+\alpha}$

VLPR09 @ Beijing, China
8/6/2009
62

**SAILING LAB**  
Laboratory for Statistical Analysis and Applied Bayesian Computation

T=1

T=2

$$\frac{2}{6 + \alpha}$$

VLPR09 @ Beijing, China
8/6/2009
63

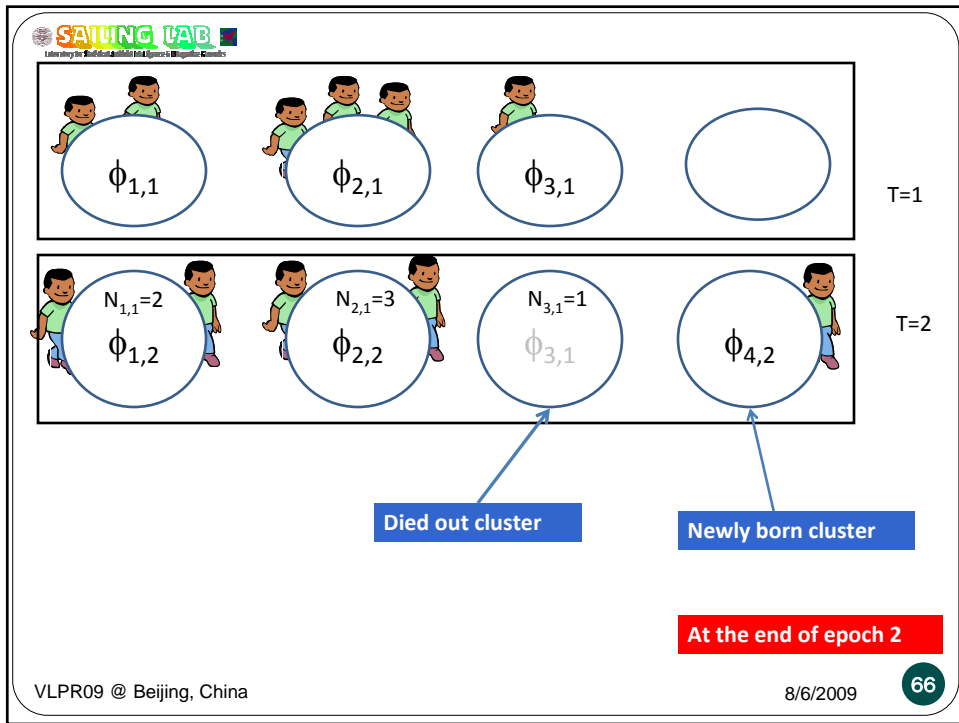
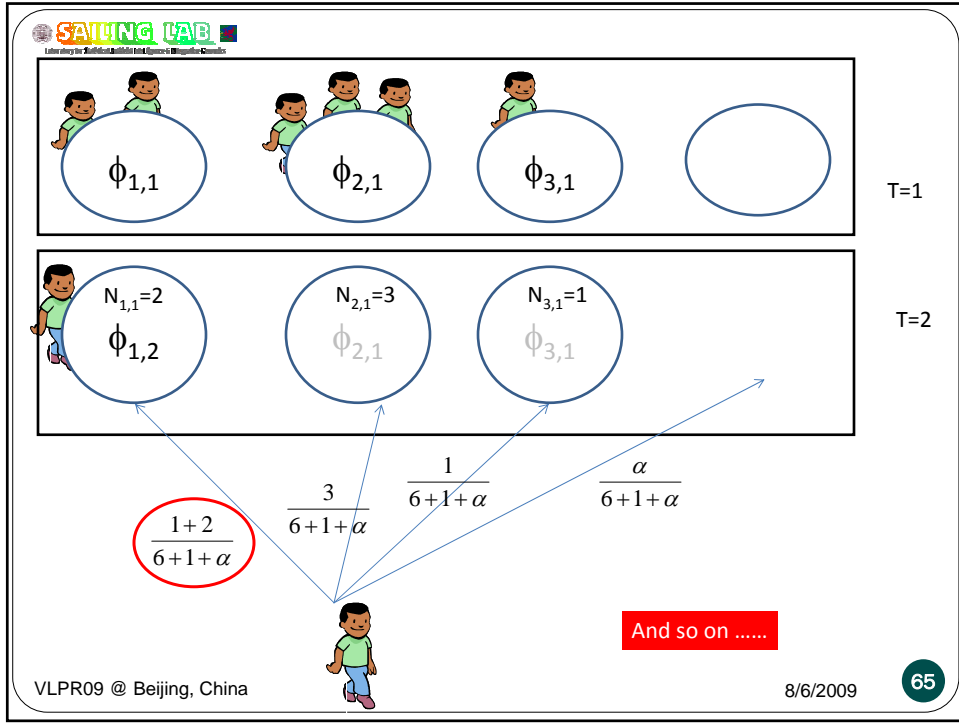
**SAILING LAB**  
Laboratory for Statistical Analysis and Applied Bayesian Computation

T=1

T=2

Sample  $\phi_{1,2} \sim P(\cdot | \phi_{1,1})$

VLPR09 @ Beijing, China
8/6/2009
64



**SAILING LAB**  
Laboratory for Statistical Analysis and Intelligent Systems Research

T=1

T=2

T=3

VLPR09 @ Beijing, China 8/6/2009 67

**SAILING LAB**  
Laboratory for Statistical Analysis and Intelligent Systems Research

## Temporal DPM

- Can be extended to model **higher-order** dependencies
- Can **decay** dependencies **over time**
  - **Pseudo-counts** for table  $k$  at time  $t$  is

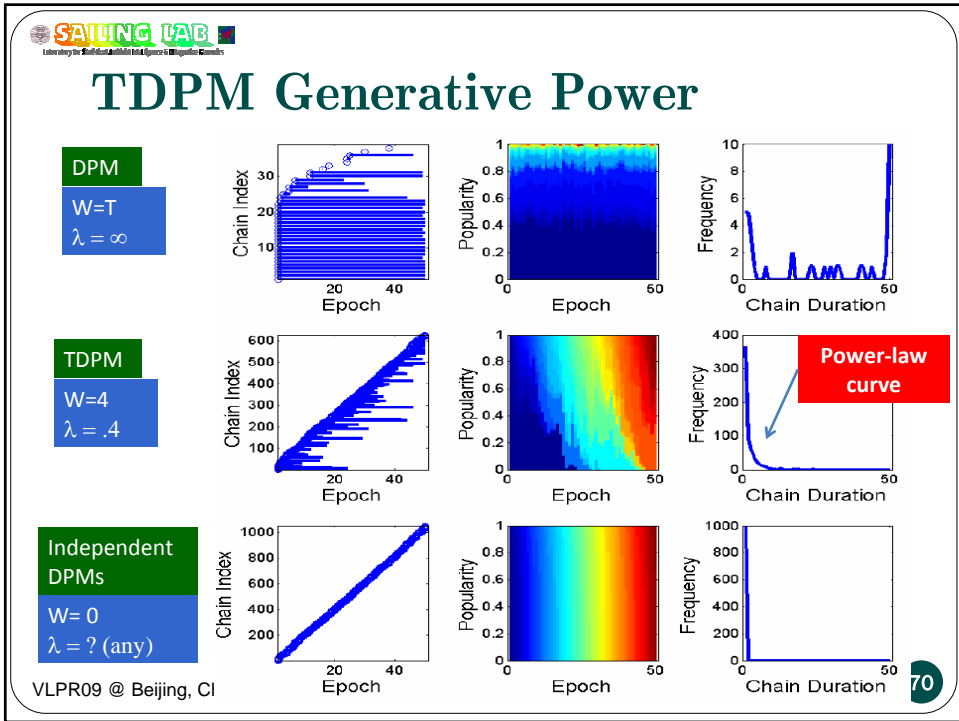
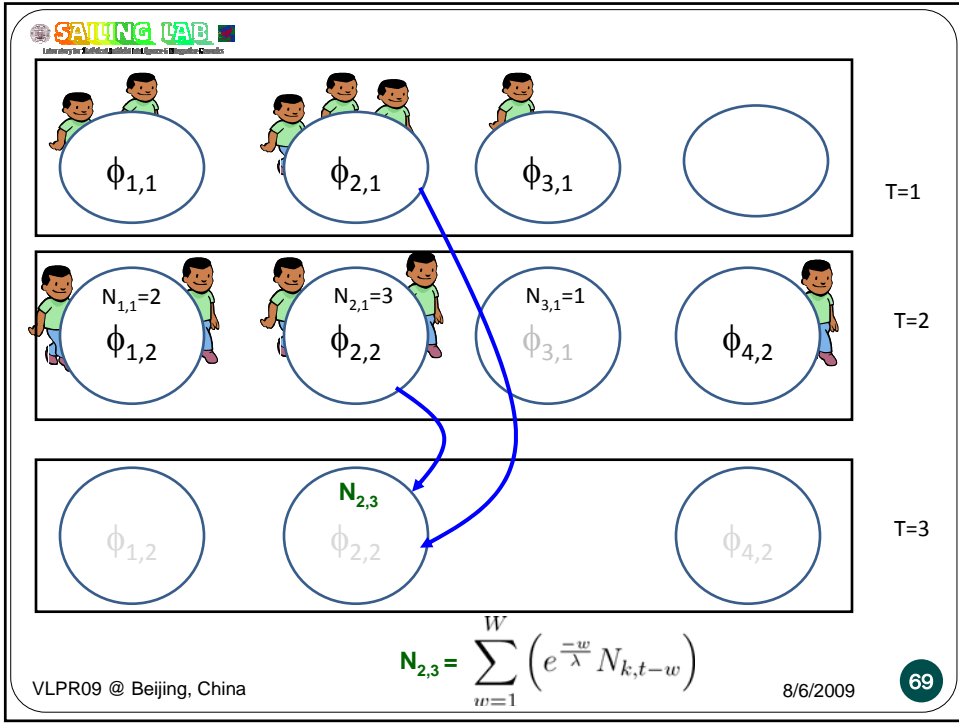
$$\sum_{w=1}^W \left( e^{\frac{-w}{\lambda}} N_{k,t-w} \right)$$

History size

Decay factory

Number of customers sitting at table  $k$  at time epoch  $t-w$

VLPR09 @ Beijing, China 8/6/2009 68

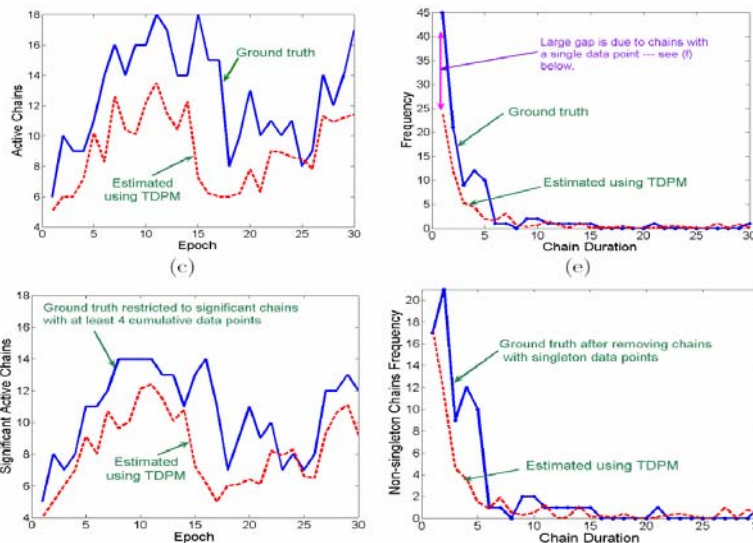


## Experiments

- **Simulated** data
- Chain dynamics is modeled as **random walk**  

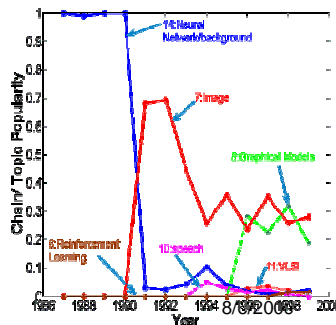
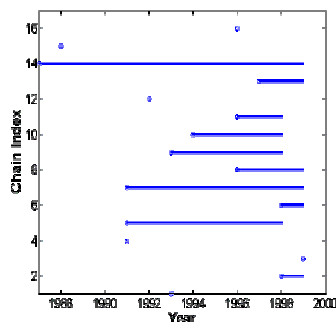
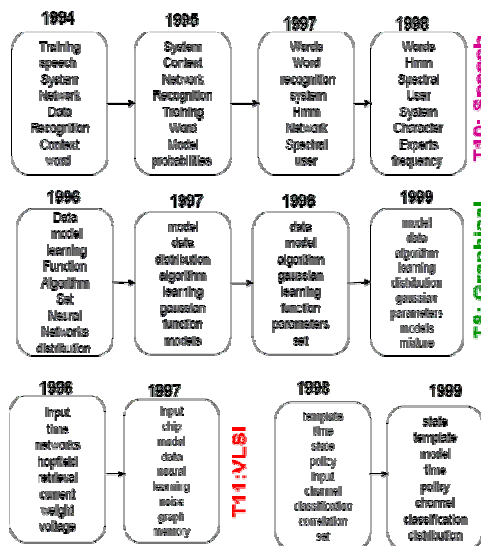
$$\phi_{k,t} | \phi_{k,t-1} \sim N(\phi_{k,t-1}, \rho I)$$
- **Gaussian emission**:  $x_{t,i} | c_{t,i} = k \sim N(\phi_{k,t}, \Sigma)$
- Simulated 30 epochs with 100 data points in each epoch
- Can TDPM recover the **ground truth** clustering?
  - **Posterior inference** ran using **Gibbs sampling** [Ahmed and Xing 2008]
- Compare with **fixed-dimension** dynamic models

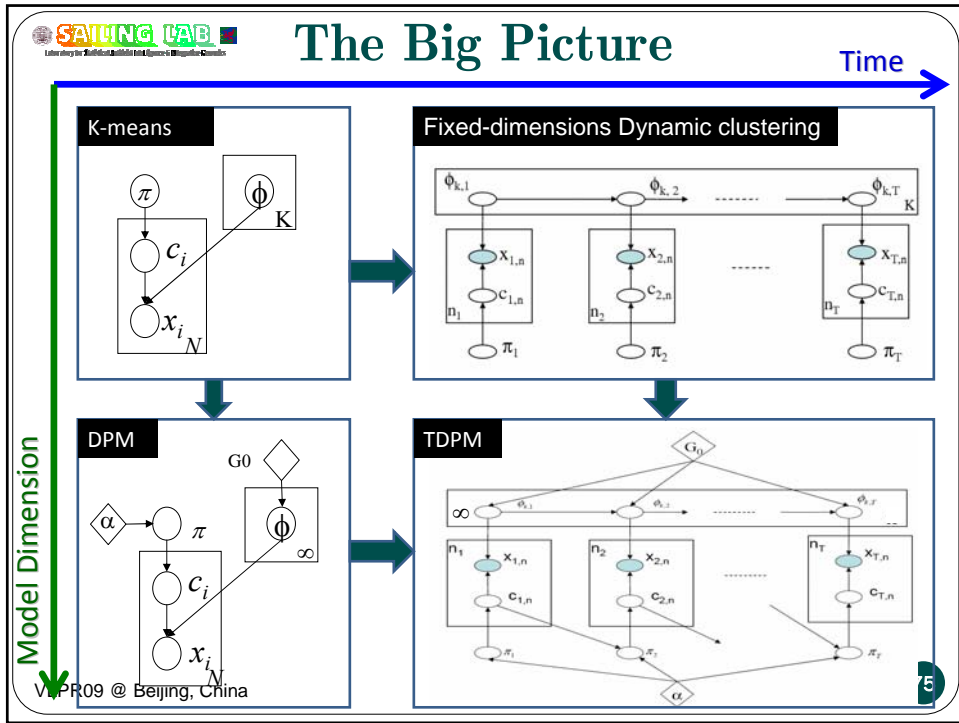
## TDPM Adaptability over Time



## Results: NIPS 12

- Building a **simple** dynamic **topic** model
- Chain dynamics is as before
- Emission model for document  $x_{k,t}$  is:
  - Project  $\phi_{k,t}$  over the **simplex**
  - Sample  $x_{k,t} | c_{k,t} \sim \text{Multinomial}(\cdot | \text{Logistic}(\phi_{k,t}))$
- Unlike LDA here a **document** belongs to **one** topic
- Use this model to analyze **NIPS12** corpus
  - Proceeding of NIPS conference 1987-1999





**SAITING LAB**  
Laboratory for Statistical Analysis and Intelligent Systems

## Summary

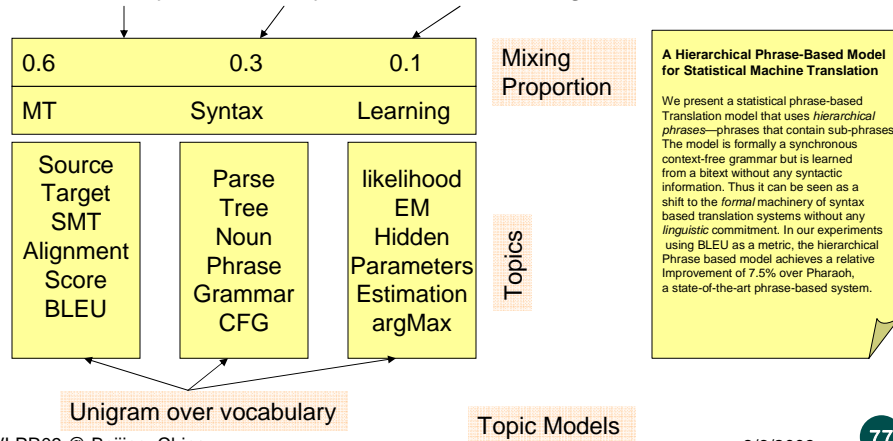
- A non-parametric Bayesian model for Pattern Uncovery
  - Finite mixture model of latent patterns (e.g., [image segments](#), [objects](#))
    - infinite mixture of propotypes: alternative to model selection
    - hierarchical infinite mixture
    - infinite hidden Markov model
    - temporal infinite mixture model
  
- Applications in general data-mining ...

VLPR09 @ Beijing, China

8/6/2009 76

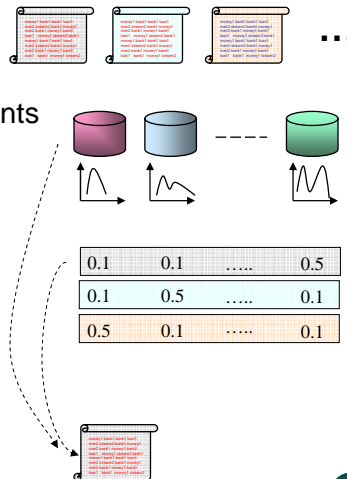
## How to Model Semantic?

- Q: What is it about?
- A: Mainly MT, with syntax, some learning



## Admixture Models

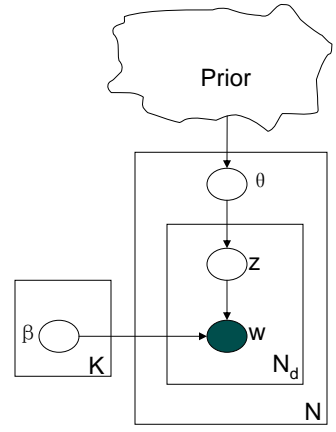
- Objects are **bags** of elements
- Mixtures are **distributions** over elements
- Objects have **mixing vector**  $\theta$ 
  - Represents each mixtures' contributions
- Object is **generated** as follows:
  - Pick a **mixture** component from  $\theta$
  - Pick an **element** from that component



# Topic Models = Admixture Models

## Generating a document

- Draw  $\theta$  from the prior
- For each word  $n$ 
  - Draw  $z_n$  from *multinomial*  $l(\theta)$
  - Draw  $w_n | z_n, \{\beta_{1:k}\}$  from *multinomial*  $l(\beta_{z_n})$



Which prior to use?

# Variational Inference

