

Reliable Variational Learning for Hierarchical Dirichlet Processes

Erik Sudderth

Brown University Computer Science

Collaborators:

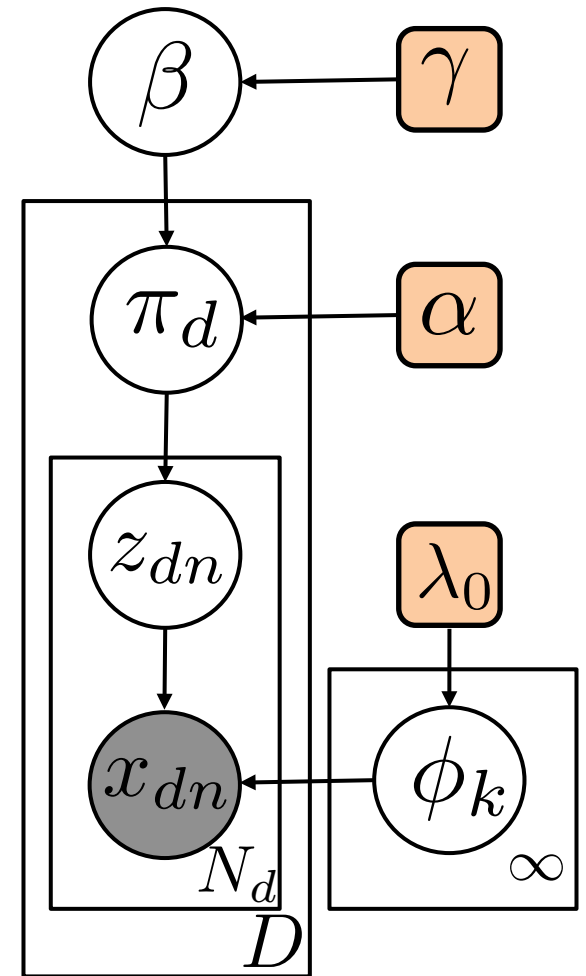
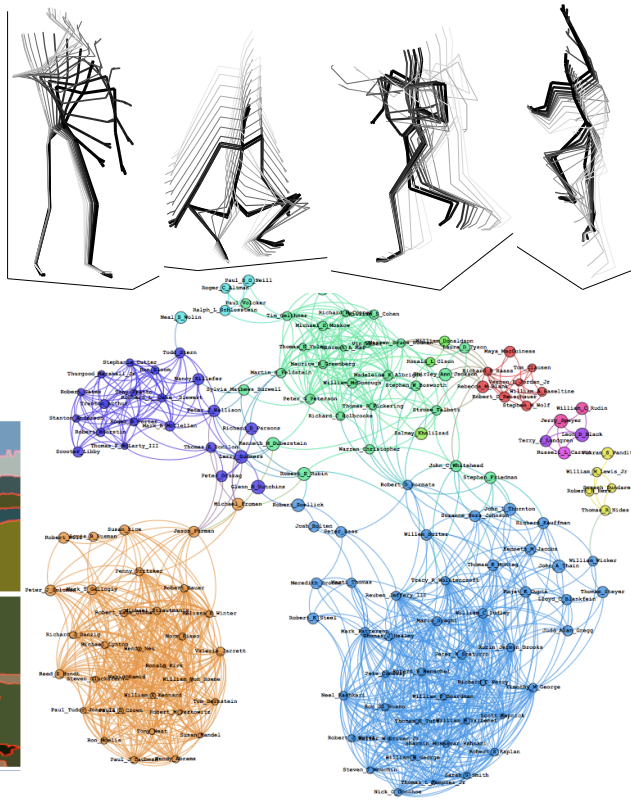
- *Michael Hughes & Dae Il Kim, Brown University*
- *Prem Gopalan & David Blei, Princeton University*



Learning Structured BNP Models

Genetics, Climate Change, Politics, ...

There are reasons to believe that the **genetics** of an **organism** are likely to shift due to the **extreme changes** in our **climate**. To protect them, our **politicians** must pass **environmental legislation** that can protect our future **species** from becoming **extinct**...



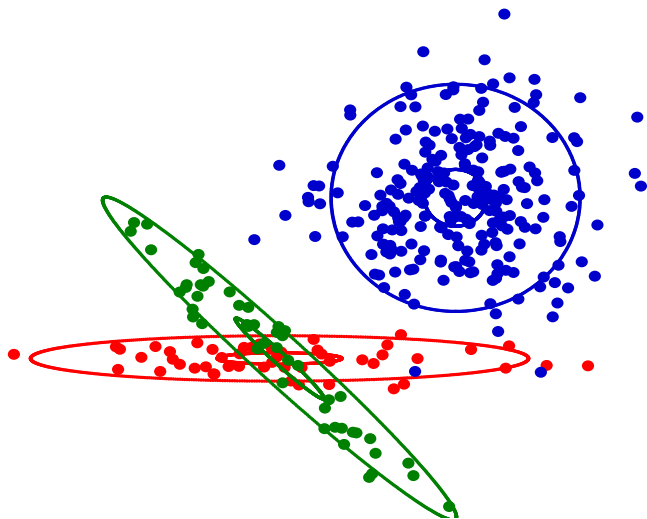
- **Nonparametric:** Data-driven discovery of model structure: *topics, behaviors, objects, communities*...
- **Reliable:** Structure driven by data and modeling assumptions, not heuristic algorithm initializations
- **Parsimonious:** Want a single model structure with good predictive power, not full posterior uncertainty

Hierarchical Dirichlet Process
(Teh et al., JASA 2006)

Memoized Variational Inference for Dirichlet Process Mixture Models

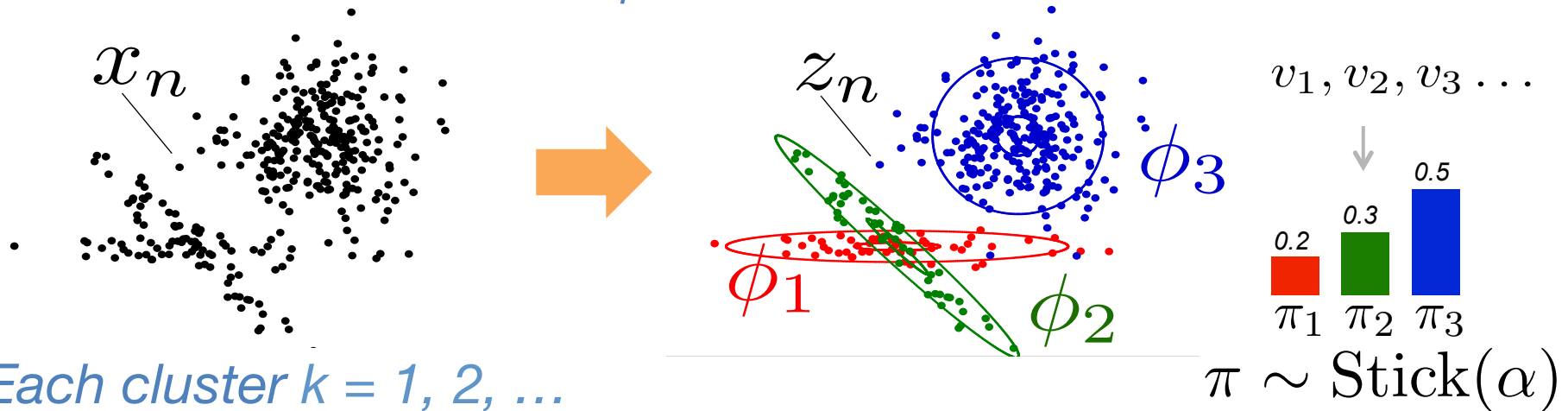
Michael Hughes & E. Sudderth

2013 Conference on Neural Information Processing Systems



Dirichlet Process Stick-Breaking

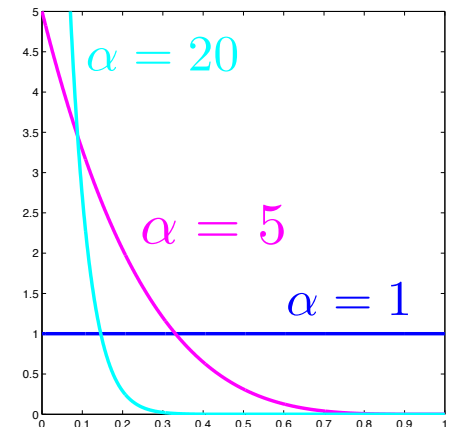
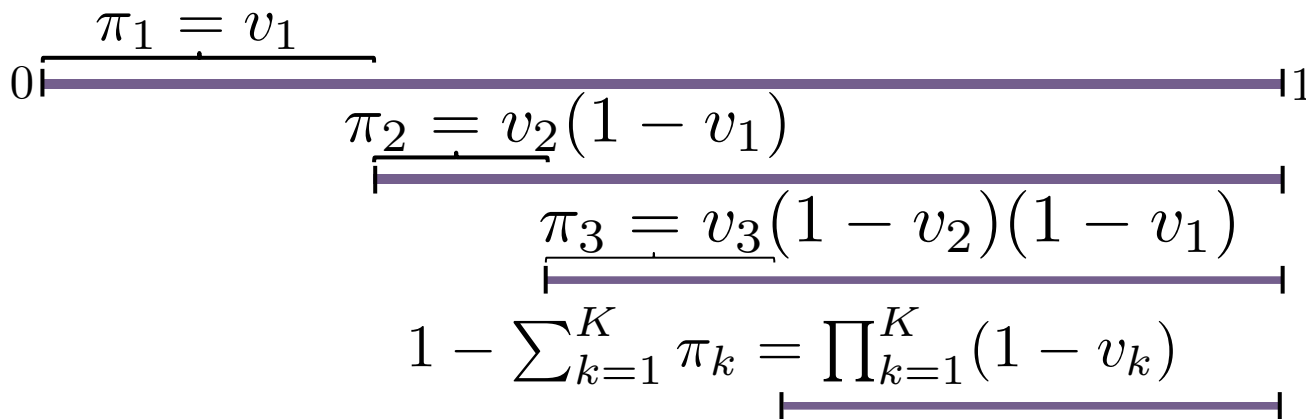
GOAL: Partition data into an a priori unknown number of discrete clusters.



Each cluster $k = 1, 2, \dots$

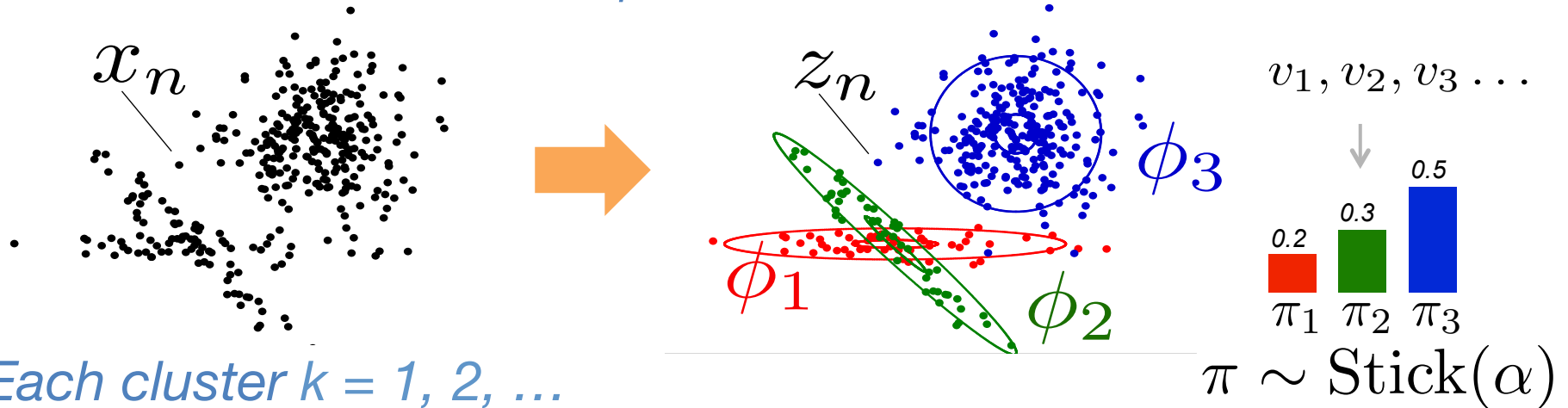
- Cluster shape: $\phi_k \sim H(\lambda_0)$
- Stick proportion: $v_k \sim \text{Beta}(1, \alpha)$
- Cluster frequency: $\pi_k = v_k \prod_{\ell=1}^{k-1} (1 - v_\ell)$

Stick-Breaking
(Sethuraman 1994)



Dirichlet Process Mixtures

GOAL: Partition data into an a priori unknown number of discrete clusters.



Each cluster $k = 1, 2, \dots$

- Cluster shape: $\phi_k \sim H(\lambda_0)$
- Stick proportion: $v_k \sim \text{Beta}(1, \alpha)$
- Cluster frequency: π_k

Each observation $n = 1, 2, \dots, N$:

- Cluster assignment: $z_n \sim \text{Cat}(\pi)$
- Observed value: $x_n \sim F(\phi_{z_n})$

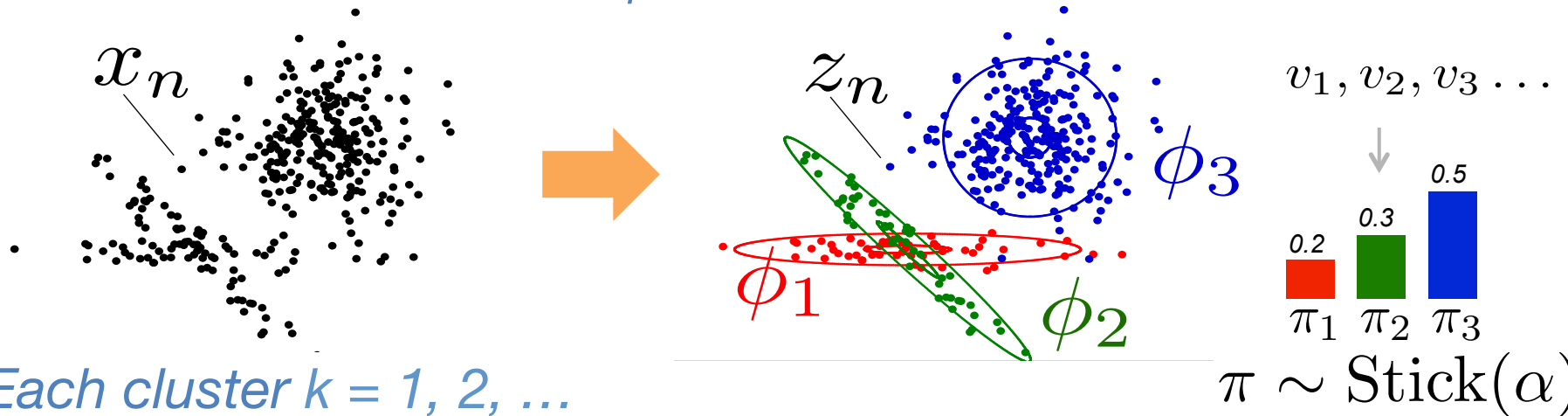
Assume exponential family likelihoods with conjugate priors

$$f(x_n | \phi_k) = \exp(\phi_k^T t(x_n) - a(\phi_k))$$

$$h(\phi_k | \lambda_0) = \exp(\lambda_0^T \bar{t}(\phi_k) - \bar{a}(\lambda_0)), \quad \bar{t}(\phi_k) = [\phi_k, -a(\phi_k)]$$

Dirichlet Process Mixtures

GOAL: Partition data into an a priori unknown number of discrete clusters.



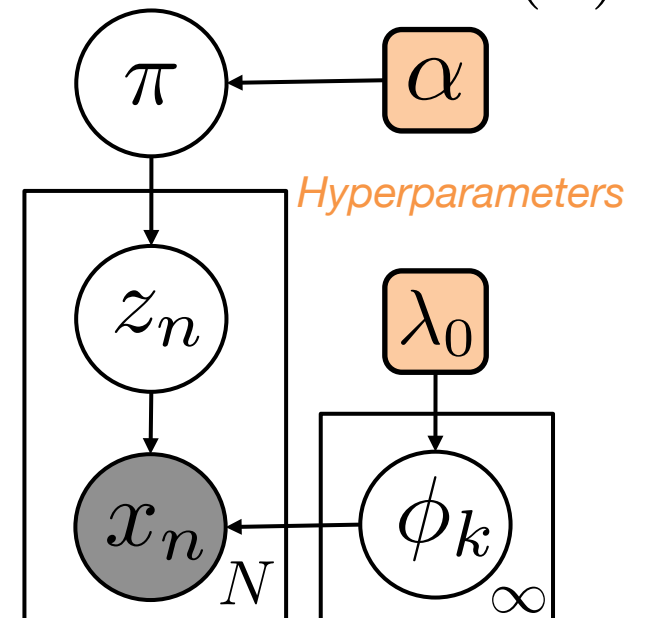
Each cluster $k = 1, 2, \dots$

- Cluster shape: $\phi_k \sim H(\lambda_0)$
- Stick proportion: $v_k \sim \text{Beta}(1, \alpha)$
- Cluster frequency: π_k

Each observation $n = 1, 2, \dots, N$:

- Cluster assignment: $z_n \sim \text{Cat}(\pi)$
- Observed value: $x_n \sim F(\phi_{z_n})$

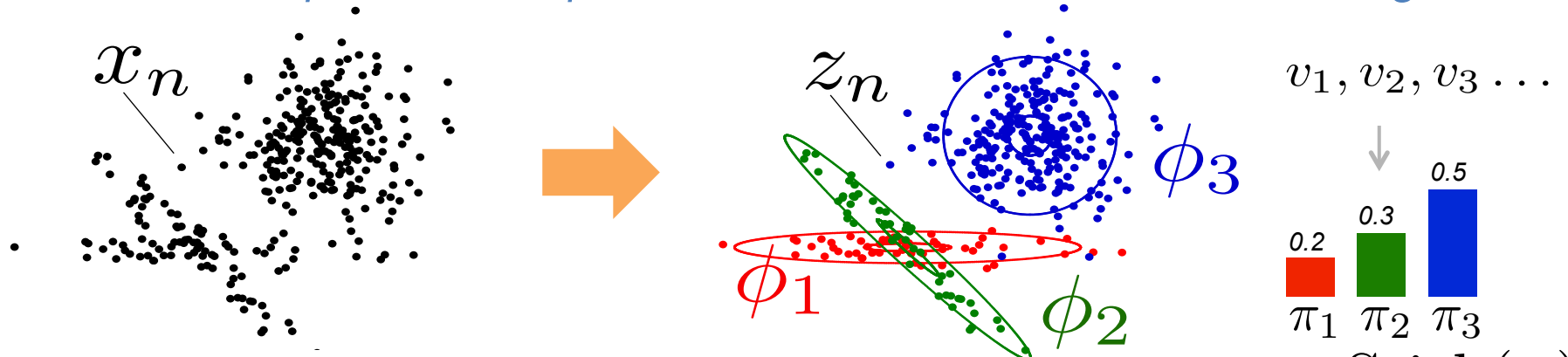
$$f(x_n | z_n = k, \phi) = \exp(\phi_k^T t(x_n) - a(\phi_k))$$



Visually summarize model structure via directed graphical model

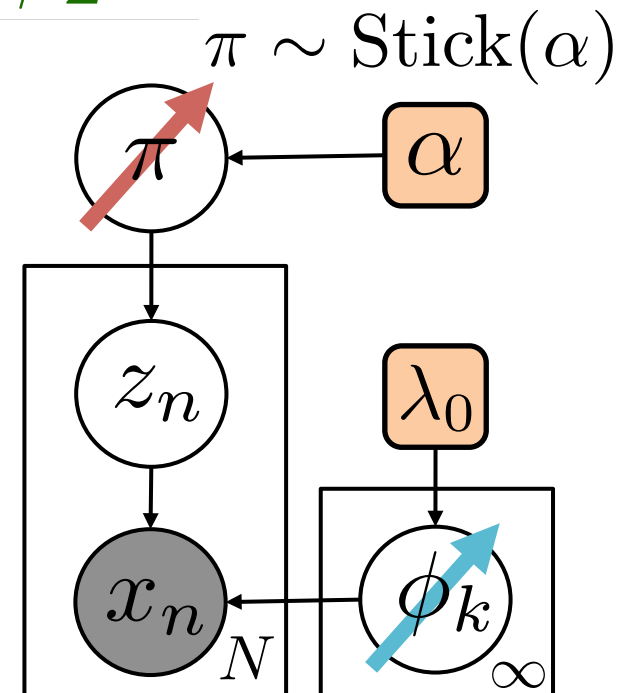
MCMC for DP Mixtures

Can we sample from the posterior distribution over data clusterings?



Given any fixed partition z :

- Marginalize stick-breaking weights via *Chinese Restaurant Process*, assigning positive probability to all partitions of data (large support)
- Via *conjugacy* of base measure to exponential family likelihood, marginalize cluster shape parameters

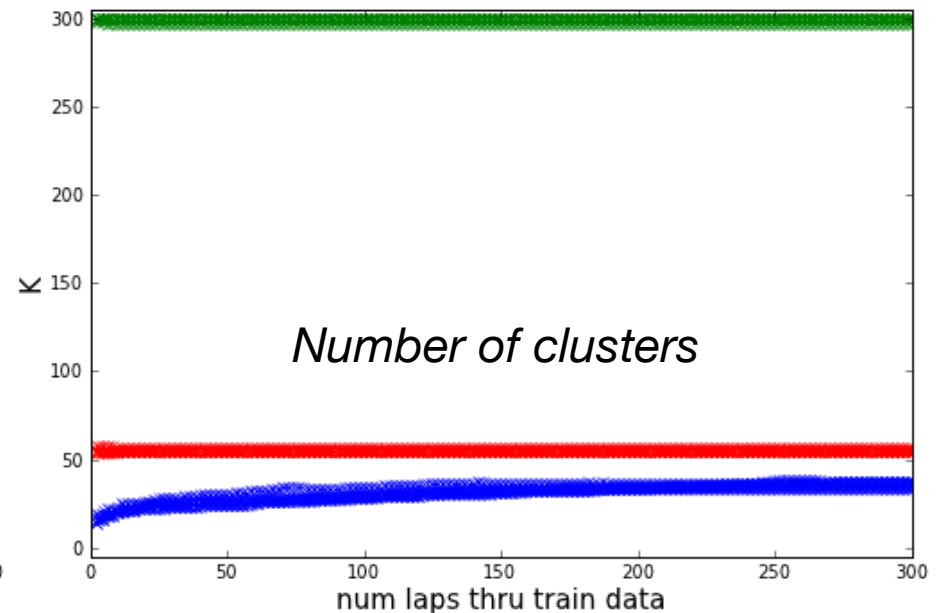
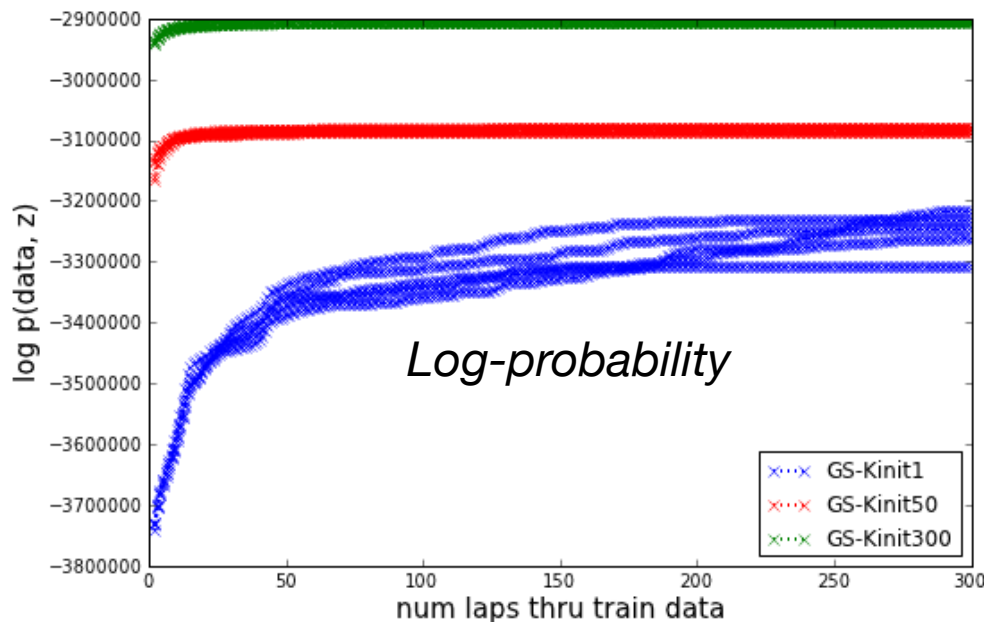
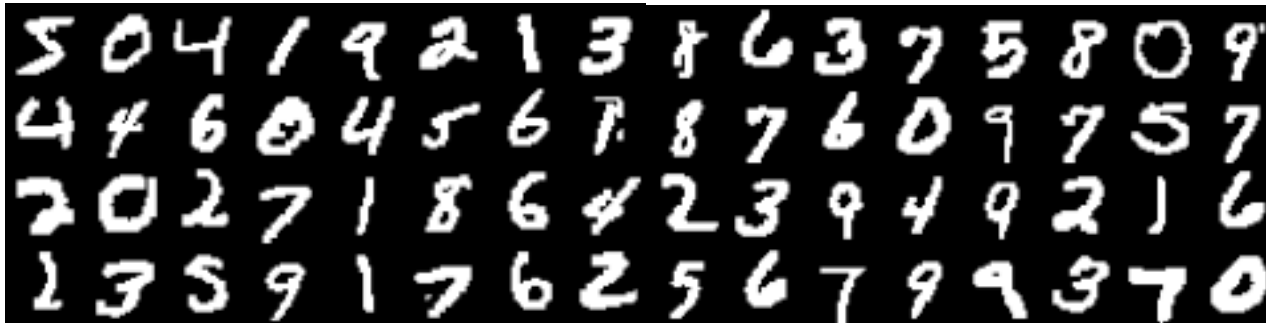


Gibbs Sampler: (Neal 1992, MacEachern 1994)

Iteratively resample cluster assignment for one observation, fixing all others.

Mixing for DP Mixture Samplers

MNIST: 60,000 digits projected to 50 dimensions via PCA.



- Five random initializations from $K=1$, $K=50$, $K=300$ clusters
- Reversible jump MCMC? Proposals slow, acceptance low.

Variational Bounds

What is the marginal likelihood of our observed data?

$$\begin{aligned}
 \log p(x \mid \alpha, \lambda_0) &= \log \sum_z \iint p(x, z, v, \phi \mid \alpha, \lambda_0) \, dv d\phi \\
 &= \log \sum_z \iint \frac{q(z, v, \phi) p(x, z, v, \phi \mid \alpha, \lambda_0)}{q(z, v, \phi)} \, dv d\phi \\
 &= \log \mathbb{E}_q \left[\frac{p(x, z, v, \phi \mid \alpha, \lambda_0)}{q(z, v, \phi)} \right]
 \end{aligned}$$

Expectation with respect to some variational distribution $q(z, v, \phi)$

Jensen's Inequality $\geq \mathbb{E}_q[\log p(x, z, v, \phi \mid \alpha, \lambda_0)] - \mathbb{E}_q[\log q(z, v, \phi)] = \mathcal{L}(q)$

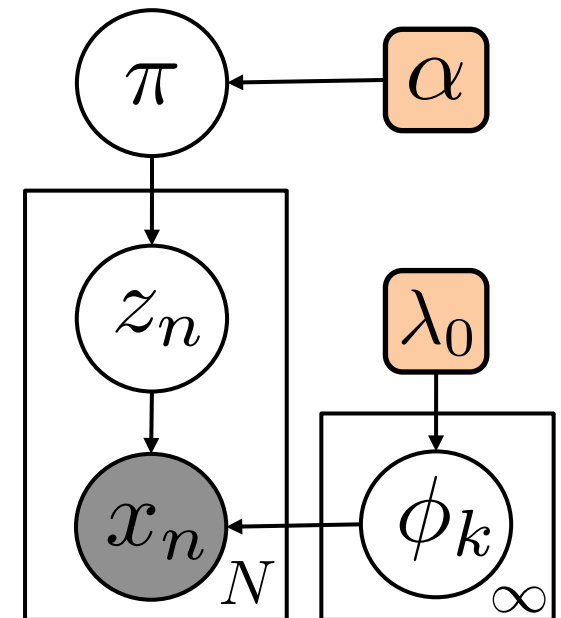
Expected log-likelihood (negative of "average energy") *Variational entropy*

- Maximizing this bound recovers true posterior:

$$\begin{aligned}
 \mathcal{L}(q) &= \log p(x \mid \alpha, \lambda_0) \\
 &\quad - \text{KL}(q(z, v, \phi) \parallel p(z, v, \phi \mid x, \alpha, \lambda_0))
 \end{aligned}$$

- The simplest *mean field* variational methods create tractable algorithms via *assumed independence*:

$$q(z, v, \phi) = q(z)q(v, \phi)$$



Approximating Infinite Models

$$q(z, v, \phi) = q(z)q(v, \phi) = \left[\prod_{n=1}^N q(z_n) \right] \cdot \left[\prod_{k=1}^{\infty} q(v_k)q(\phi_k) \right]$$

$q(z_n = k) = r_{nk}$
Beta Distribution
Exponential Family from Conjugate Prior

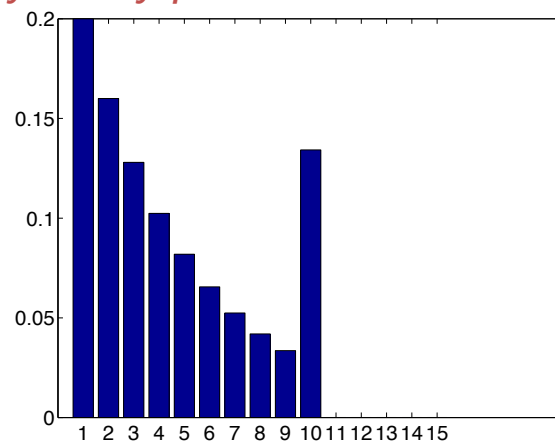
Categorical distribution with unbounded support, and infinitely many potential clusters!

Top-Down Model Truncation

Blei & Jordan, 2006; Ishwaran & James, 2001

$$q(z_n) = \text{Cat}(z_n \mid r_{n1}, r_{n2}, \dots, r_{nK})$$

$$q(v, \phi) = \left[\prod_{k=1}^K q(\phi_k) \right] \cdot \left[\prod_{k=1}^{K-1} q(v_k) \right], \quad v_K = \prod_{k=1}^{K-1} (1 - v_k).$$



$\alpha = 4, K = 10$

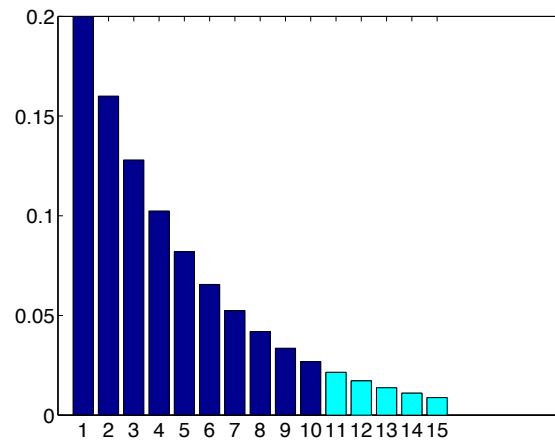
Bottom-Up Assignment Truncation

Bryant & Sudderth, 2012; Teh, Kurihara, & Welling, 2008

$$q(z_n) = \text{Cat}(z_n \mid r_{n1}, r_{n2}, \dots, r_{nK}, 0, 0, 0, \dots)$$

$$q(v, \phi) = \prod_{k=1}^{\infty} q(v_k)q(\phi_k)$$

For any $k > K$, optimal variational distributions equal prior & need not be explicitly represented



Batch Variational Updates

A Bayesian nonparametric analog of Expectation-Maximization (EM)

$$q(z, v, \phi) = \left[\prod_{n=1}^N q(z_n | r_n) \right] \cdot \left[\prod_{k=1}^{\infty} \text{Beta}(v_k | \alpha_{k1}, \alpha_{k0}) h(\phi_k | \lambda_k) \right]$$

$$q(z_n) = \text{Cat}(z_n | r_{n1}, r_{n2}, \dots, r_{nK}, 0, 0, 0, \dots) \quad \text{for some } K > 0$$

Update Assignments (The Expectation Step): For all N data,

$$r_{nk} \propto \exp(\mathbb{E}_q[\log \pi_k(v)] + \mathbb{E}_q[\log p(x_n | \phi_k)]) \quad \text{for } k \leq K$$

$$\mathbb{E}_q[\log \pi_k(v)] = \underbrace{\mathbb{E}_q[\log(v_k)]}_{\psi(\alpha_{k1}) - \psi(\alpha_{k1} + \alpha_{k0})} + \sum_{\ell=1}^{k-1} \underbrace{\mathbb{E}_q[\log(1 - v_\ell)]}_{\psi(\alpha_{k0}) - \psi(\alpha_{k1} + \alpha_{k0})}$$

Update Cluster Parameters (The Other Expectation Step):

$$N_k^0 = \sum_{n=1}^N r_{nk} \quad s_k^0 \leftarrow \sum_{n=1}^N r_{nk} t(x_n) \quad \lambda_k \leftarrow \lambda_0 + s_k^0$$

Expected counts and sufficient statistics are only non-zero for first K clusters

$$\alpha_{k1} \leftarrow 1 + N_k^0 \quad \mathbb{E}_q[v_k] = \frac{\alpha_{k1}}{\alpha_{k1} + \alpha_{k0}}$$
$$\alpha_{k0} \leftarrow \alpha + \sum_{\ell=k+1}^{\infty} N_\ell^0 = \alpha + \sum_{\ell=k+1}^K N_\ell^0$$

Likelihood Bounds & Convergence

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(x, z, v, \phi \mid \alpha, \lambda_0)] - \mathbb{E}_q[\log q(z, v, \phi)]$$

➤ Immediately after global parameter update, bound simplifies:

$$\mathcal{L}(q) = \mathbb{H}[r] + \sum_{k=1}^K [\bar{a}(\lambda_k) - \bar{a}(\lambda_0) + \log B(\alpha_{k1}, \alpha_{k0}) - \log B(1, \alpha)]$$

log-normalizers for cluster shape and beta stick-breaking priors

$$\mathbb{H}[r] = - \sum_{n=1}^N \sum_{k=1}^{\infty} r_{nk} \log r_{nk} = - \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log r_{nk}$$

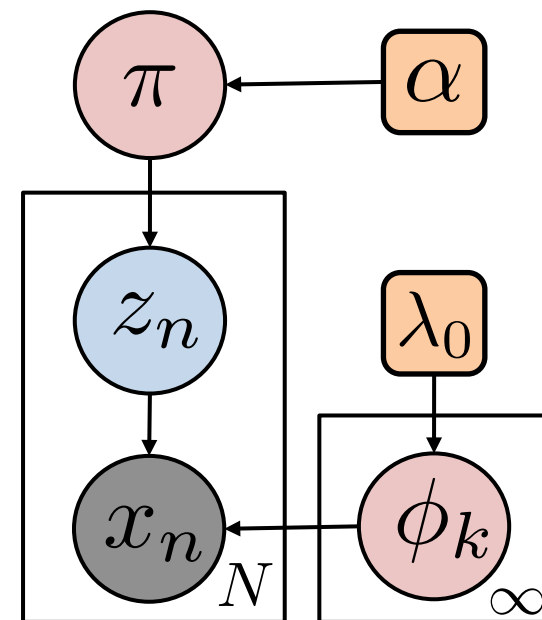
For data item $n = 1, 2, \dots, N$, and K candidate clusters:

$$q(z_n = k) = r_{nk} \propto e^{\mathbb{E}_q[\log \pi_k(v) + \log p(x_n | \phi_k)]}$$

For cluster $k = 1, 2, \dots, K$:

$$s_k^0 \leftarrow \sum_{n=1}^N r_{nk} t(x_n) \quad \text{Match Expected Sufficient Statistics} \quad \alpha_{k1} \leftarrow 1 + N_k^0$$

$$\lambda_k \leftarrow \lambda_0 + s_k^0 \quad \alpha_{k0} \leftarrow \alpha + N_{>k}^0$$



Likelihood Bounds & Convergence

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(x, z, v, \phi \mid \alpha, \lambda_0)] - \mathbb{E}_q[\log q(z, v, \phi)]$$

- Immediately after global parameter update, bound simplifies:

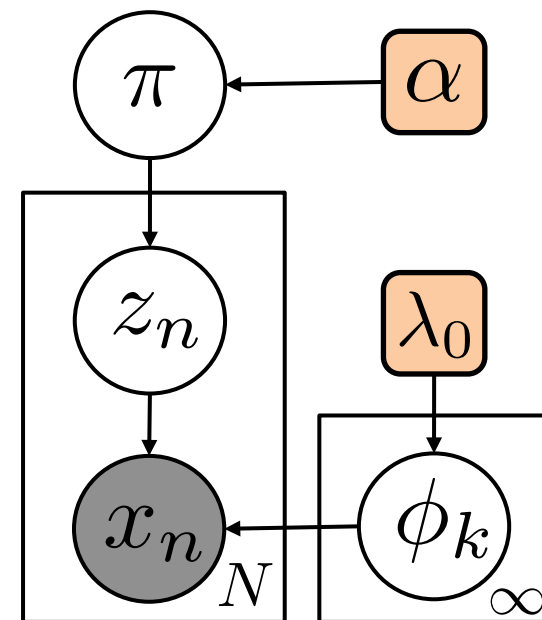
$$\mathcal{L}(q) = \mathbb{H}[r] + \sum_{k=1}^K [\bar{a}(\lambda_k) - \bar{a}(\lambda_0) + \log B(\alpha_{k1}, \alpha_{k0}) - \log B(1, \alpha)]$$

log-normalizers for cluster shape and beta stick-breaking priors

$$\mathbb{H}[r] = - \sum_{n=1}^N \sum_{k=1}^{\infty} r_{nk} \log r_{nk} = - \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log r_{nk}$$

- Properties of variational optimization algorithm:

- + Likelihood bound monotonically increasing, guaranteed convergence to posterior mode
- + Unlike classical EM for MAP estimation, allows Bayesian comparison of hypotheses with varying complexity K , crucial for BNP models
- Truncation level K is assumed fixed
- Sensitive to initialization (many modes)
- Each iteration must examine all data (SLOW)



Stochastic Variational Inference

Hoffman, Blei, Paisley, & Wang, JMLR 2013

Stochastically partition large dataset into B smaller *batches*:

Update: For each batch b

$$r(\mathcal{B}_b) \leftarrow \text{Estep}(x(\mathcal{B}_b), \alpha, \lambda)$$

For cluster $k = 1, 2, \dots, K$:

$$s_k^b \leftarrow \sum_{n \in \mathcal{B}_b} r_{nk} t(x_n) \quad \text{batch stats give noisy estimate of}$$

$$\lambda_k^b \leftarrow \lambda_0 + \frac{N}{|\mathcal{B}_b|} s_k^b \quad \text{(natural)}$$

$$\lambda_k \leftarrow \rho_t \lambda_k^b + (1 - \rho_t) \lambda_k \quad \text{gradient}$$

Apply similar updates to stick weights.

Data

$x(\mathcal{B}_1)$

$x(\mathcal{B}_2)$

\vdots

$x(\mathcal{B}_b)$

\vdots

$x(\mathcal{B}_B)$

Learning Rate

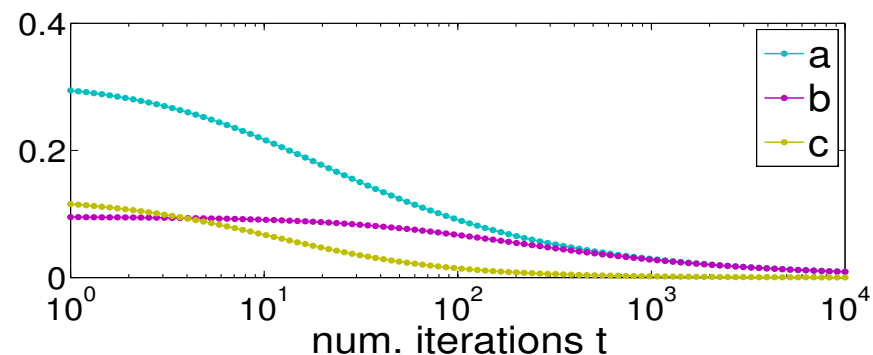
$$\rho_t \triangleq (\rho_0 + t)^{-\kappa}$$

Robbins-Monro convergence condition:

$$\sum_t \rho_t \rightarrow \infty \quad \kappa \in (.5, 1]$$
$$\sum_t \rho_t^2 < \infty$$

Properties of stochastic inference:

- + Per-iteration cost is low
- + Initial iterations often very effective
- Objective is highly non-convex, so convergence guarantee is weak
- Batch size and learning rate significantly impact efficiency & accuracy



Memoized Variational Inference

Hughes & Sudderth, NIPS 2013; Neal & Hinton 1999

Memoization: Storage (caching) of results of previous computations

Update: For each batch b

$$r(\mathcal{B}_b) \leftarrow \text{Estep}(x(\mathcal{B}_b), \alpha, \lambda)$$

For cluster $k = 1, 2, \dots, K$:

$$s_k^0 \leftarrow s_k^0 - s_k^b$$

$$s_k^b \leftarrow \sum_{n \in \mathcal{B}_b} r_{nk} t(x_n)$$

$$s_k^0 \leftarrow s_k^0 + s_k^b$$

$$\lambda_k \leftarrow \lambda_0 + s_k^0$$

batch stats allow exact estimation from partial E-steps

Apply similar updates to stick weights.

Data	Batch Summaries
$x(\mathcal{B}_1)$	$s_1^1 \quad s_2^1 \quad \dots \quad s_K^1$
$x(\mathcal{B}_2)$	$s_1^2 \quad s_2^2 \quad \dots \quad s_K^2$
\vdots	\vdots
$x(\mathcal{B}_b)$	\vdots
\vdots	\vdots
$x(\mathcal{B}_B)$	$s_1^B \quad s_2^B \quad \dots \quad s_K^B$

Properties of memoized inference:

- + Per-iteration cost is low
- + Initial iterations often very effective
- + Insensitive to chosen B , no learning rate
- + Foundation for inferring number of clusters K
- Requires storage proportional to number of batches (NOT number of observations)

Global Summary

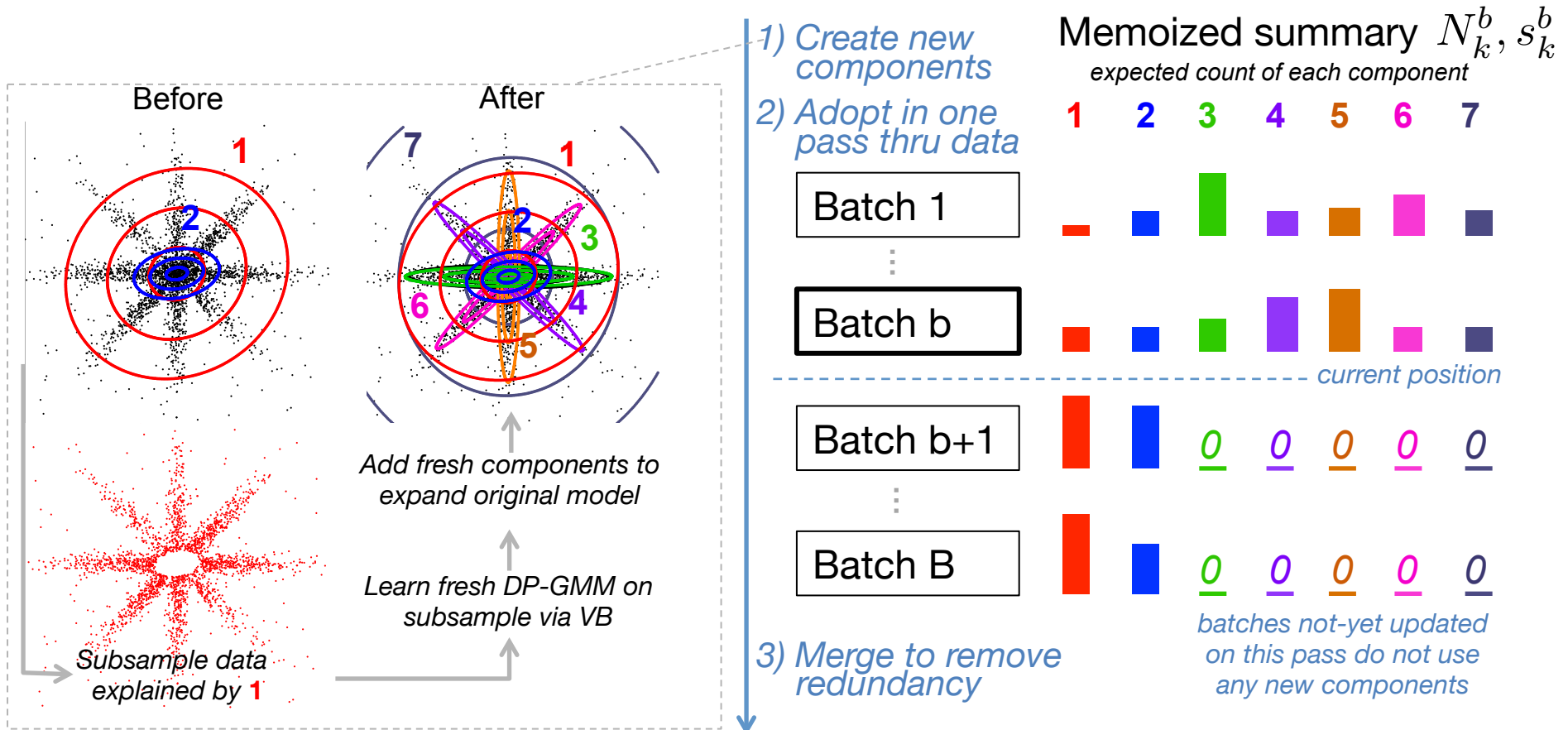
$$s_k^0 = s_k^1 + s_k^2 + \dots + s_k^B$$

Entropy for $\mathcal{L}(q)$

$$H_k^0 = H_k^1 + H_k^2 + \dots + H_k^B$$

$$H_k^b = - \sum_{n \in \mathcal{B}_b} r_{nk} \log r_{nk}$$

Memoized Cluster Births

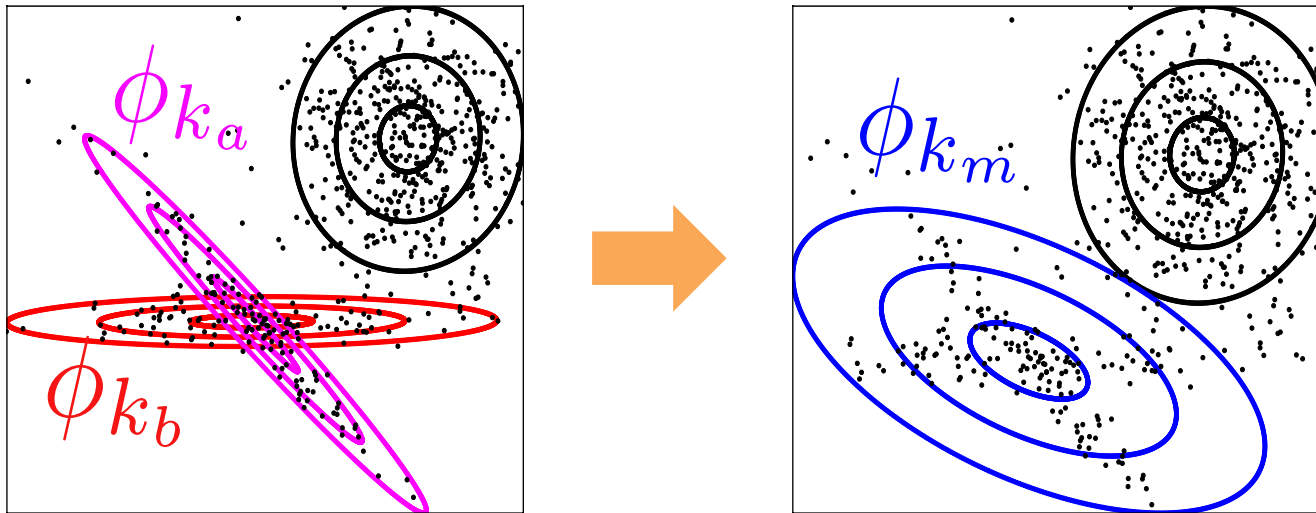


Principles guiding memoized births:

- BNP models support rare clusters, so random sampling ineffective
- Target data grouped by some current cluster (likelihood-independent)
- Memoized updates allow efficient marginal likelihood verification

Memoized Cluster Merges

Merge two clusters into one for parsimony, accuracy, efficiency.



- New cluster takes over all responsibility for data assigned to old clusters:

$$r_{nk_m} \leftarrow r_{nk_a} + r_{nk_b} \quad \longrightarrow \quad N_{k_m}^0 \leftarrow N_{k_a}^0 + N_{k_b}^0, \quad s_{k_m}^0 \leftarrow s_{k_a}^0 + s_{k_b}^0$$

- No batch processing required, efficiently evaluate via *memoized* statistics:

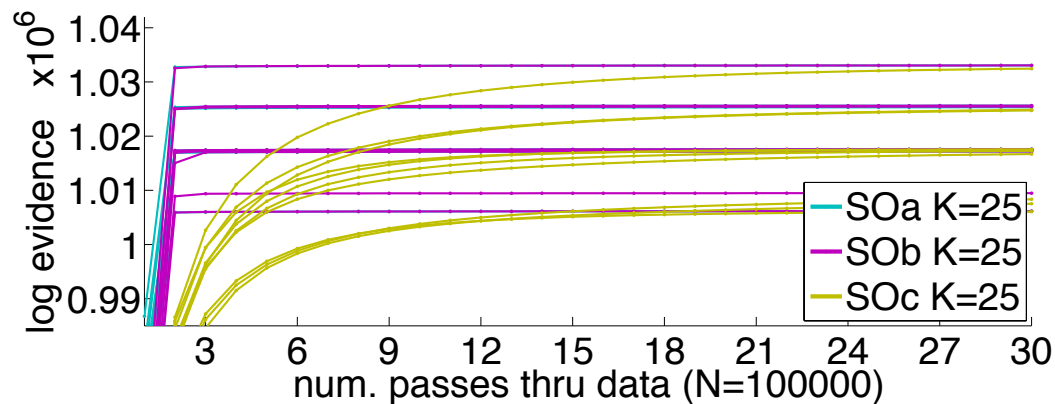
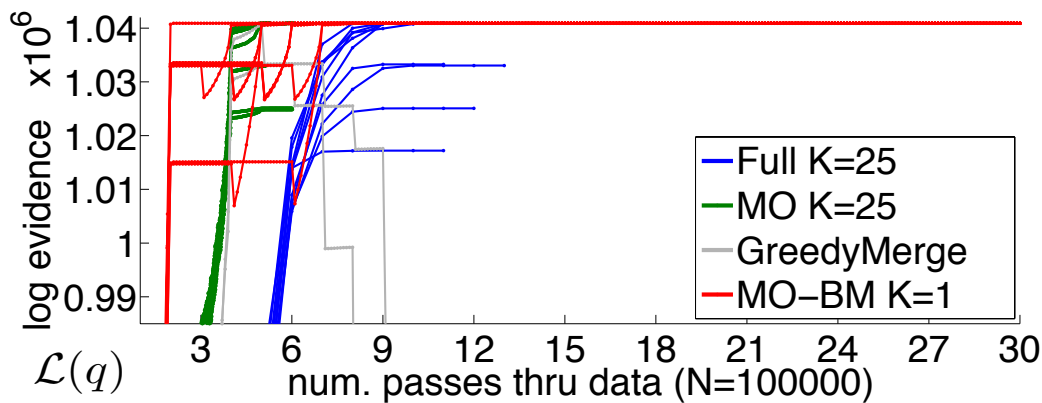
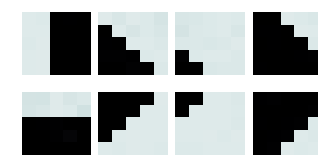
$$\mathcal{L}(q) = \mathbb{H}[r] + \sum_{k=1}^K \left[\bar{a}(s_k^0 + \lambda_0) - \bar{a}(\lambda_0) + \log B(1 + N_k^0, \alpha + N_{>k}^0) - \log B(1, \alpha) \right]$$

- Accept or reject via *exact* full-dataset likelihood bound: $\mathcal{L}(q_{\text{merge}}) > \mathcal{L}(q)$?

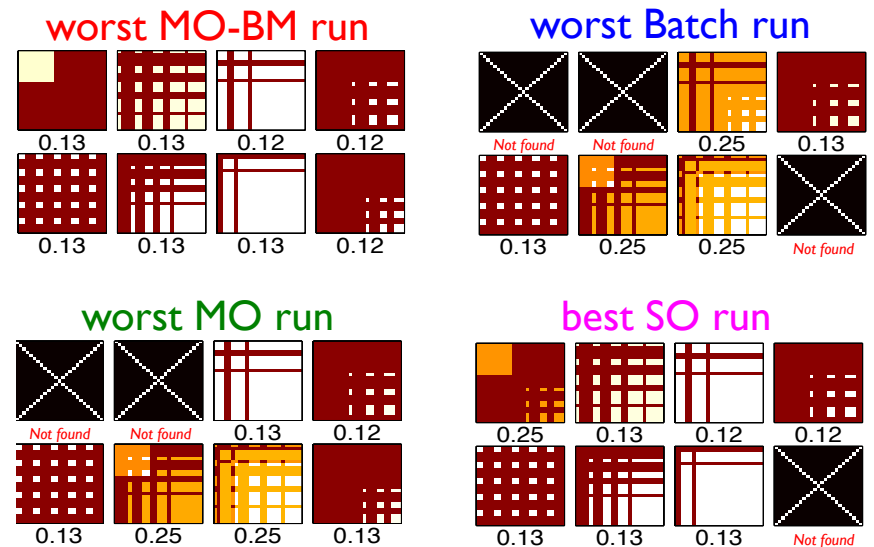
*Requires memoized entropy sums for candidate pairs of clusters;
more efficient alternatives under development.*

Example: Finite Gaussian Mixture

- $N=100,000$ samples from mixture of 8 Gaussians
- 25-dim. covariance motivated by 5x5 image patches
- DP mixture variational approximations allow $K=25$ clusters



From 10 random initializations:

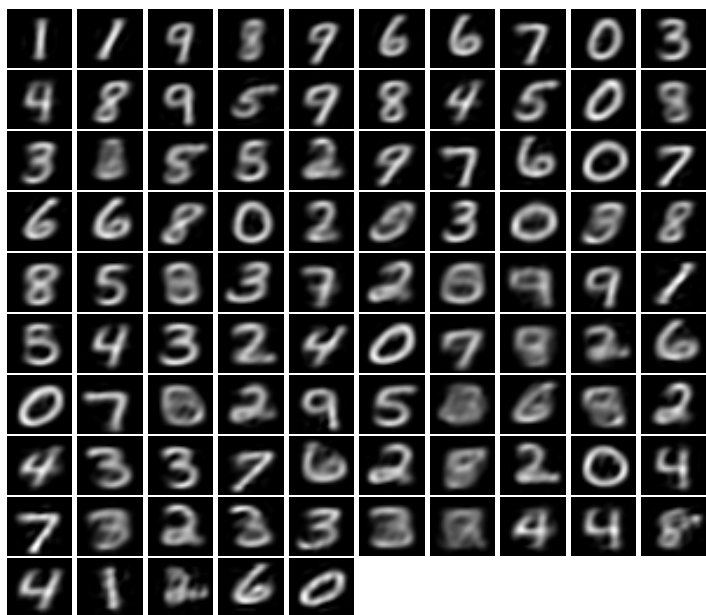


- Memoized birth-merge from $K=1$ finds true cluster every time
- Greedy decisions to merge based on single batches collapse model
- Stochastic sensitive to learning rate and initialization

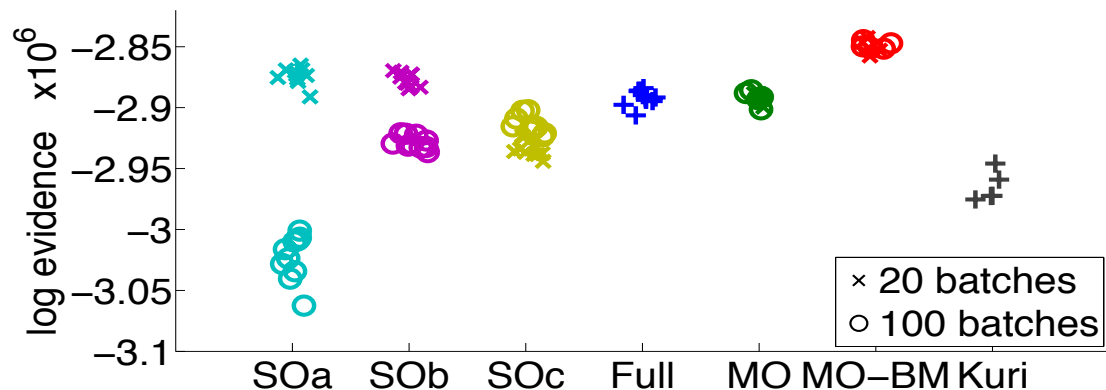
Batch, memoized, & memoized birth-merge
 Stochastic variational: Rate a, Rate b, Rate c
 Greedy: Merge based on single batches

Clustering Handwritten Digits

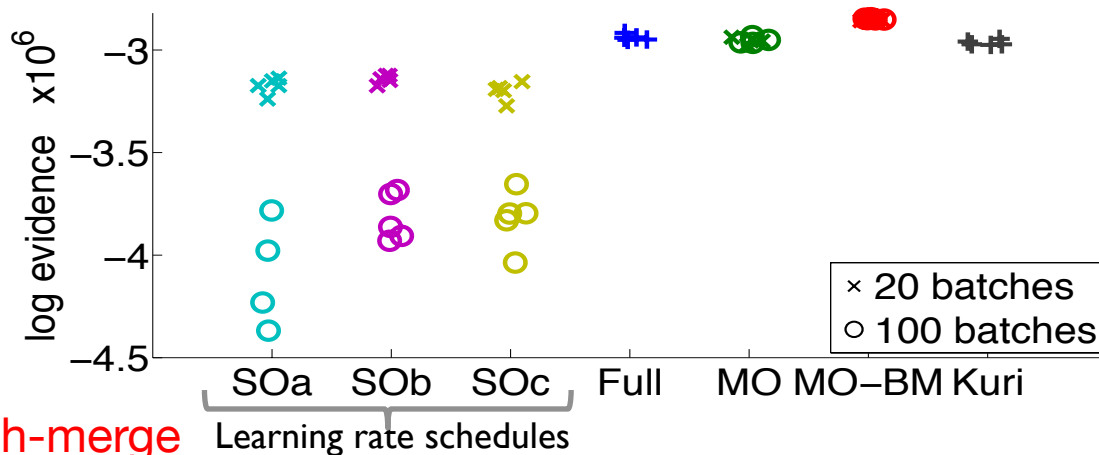
MNIST: 60,000 digits projected to 50 dimensions via PCA.



Likelihood bound, K-means++ initialization



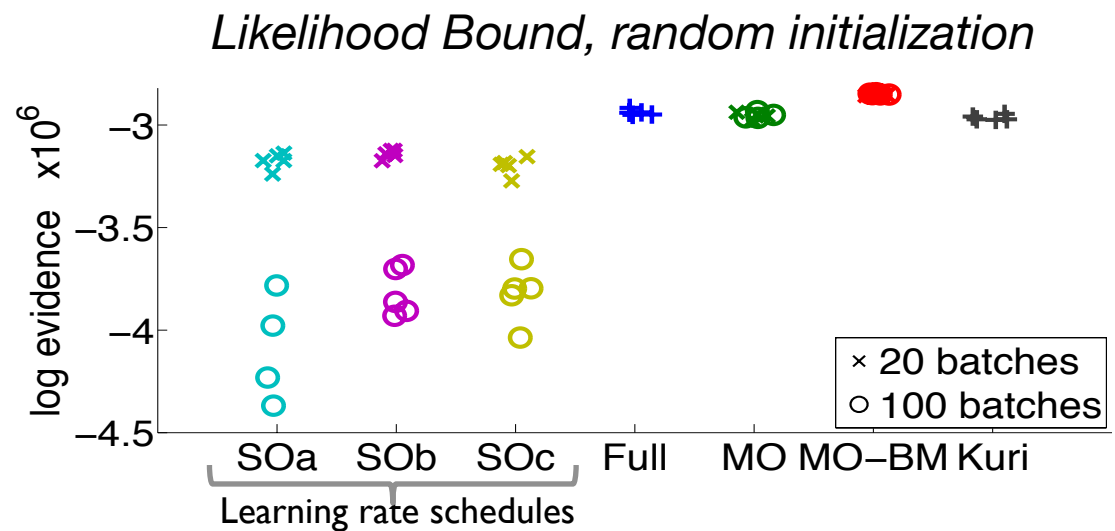
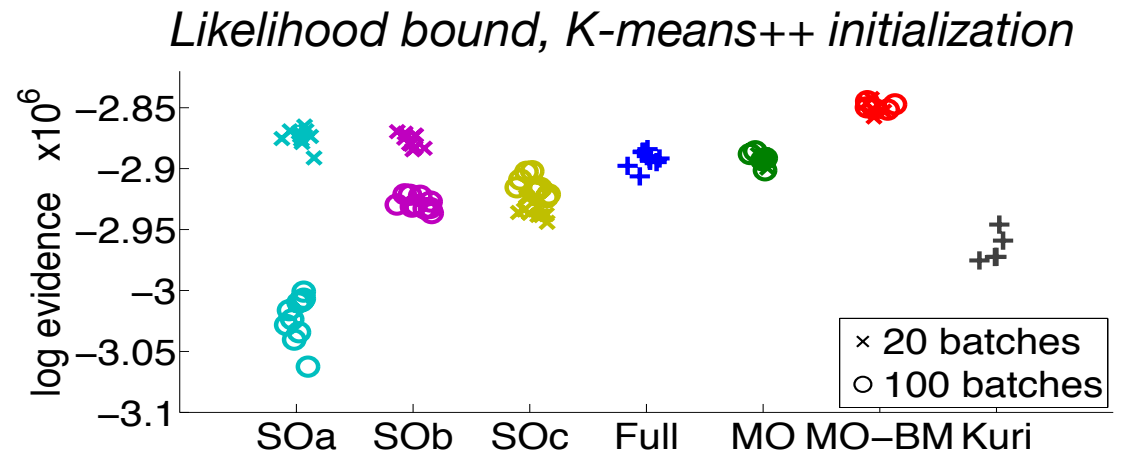
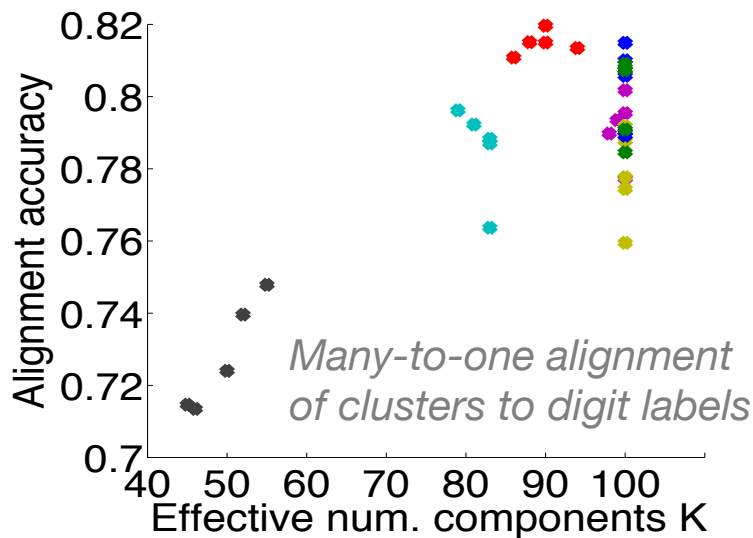
Likelihood Bound, random initialization



Batch, memoized, & memoized birth-merge
 Stochastic variational: Rate a, Rate b, Rate c
 Kurihara: Accelerated variational, NIPS 2006

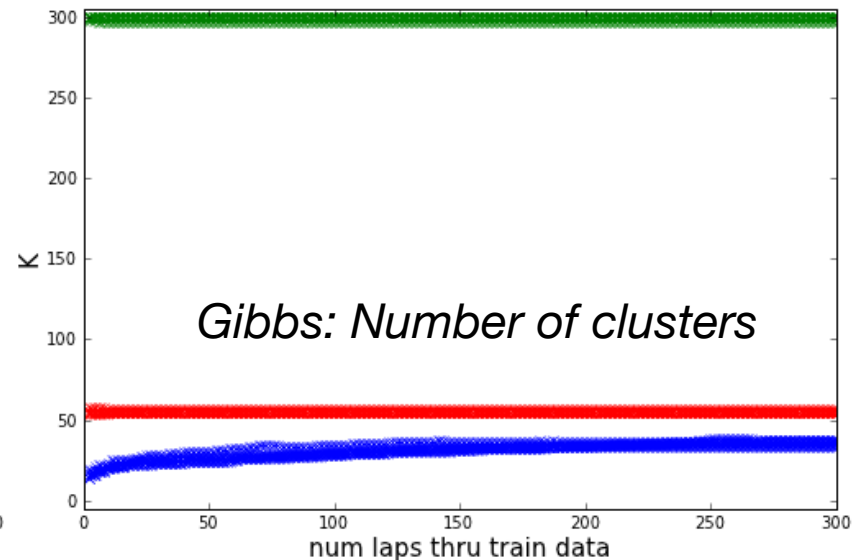
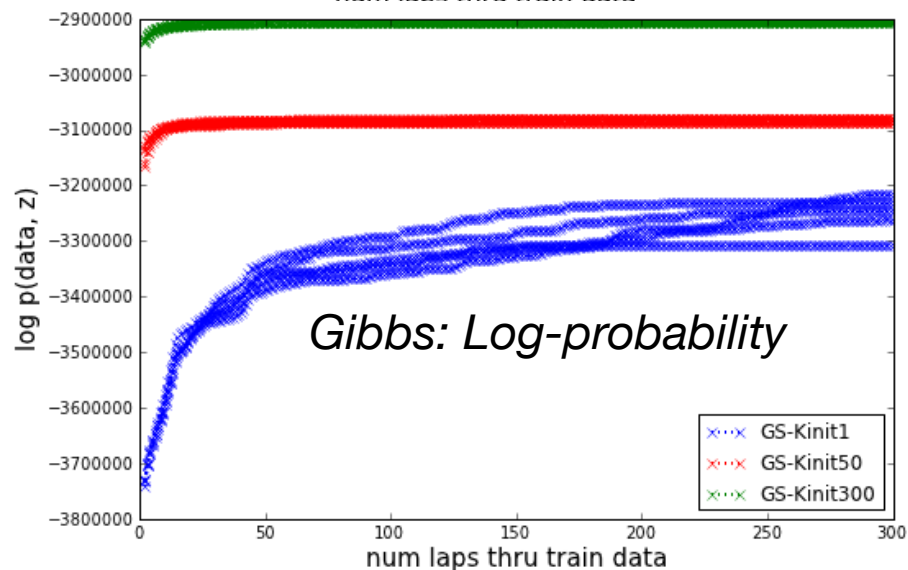
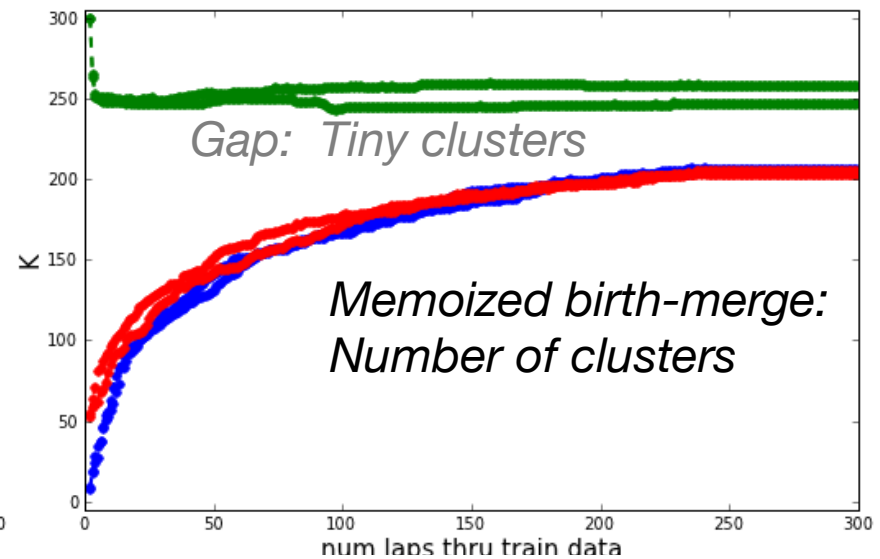
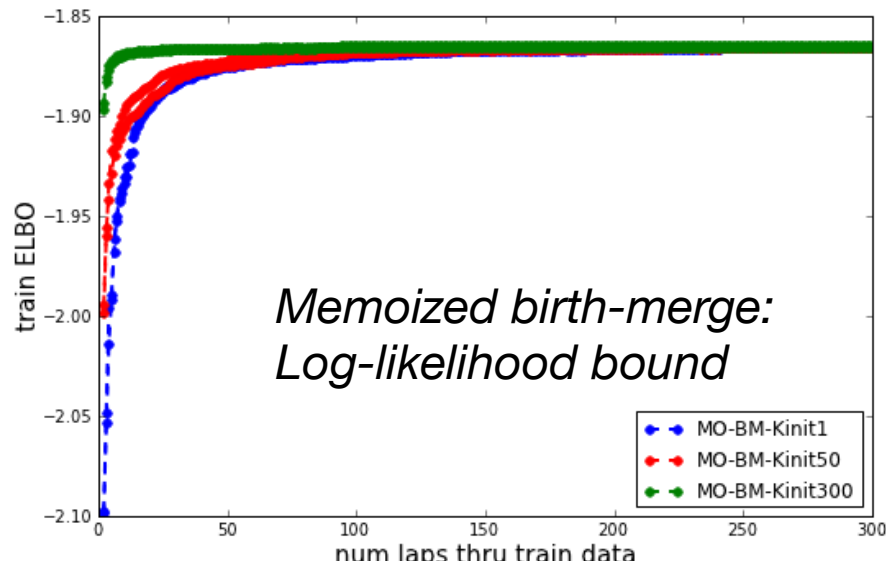
Clustering Handwritten Digits

MNIST: 60,000 digits projected to 50 dimensions via PCA.



➤ **Memoized birth-merge** from K=1 has highest accuracy while using fewer clusters

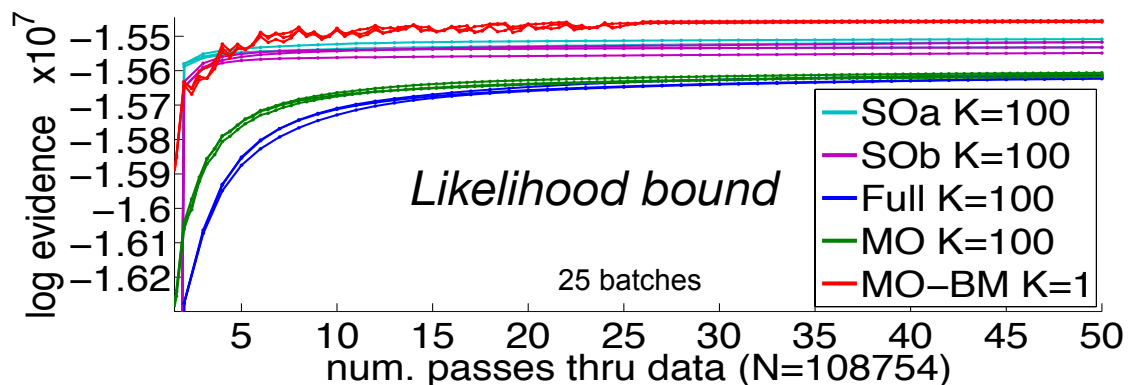
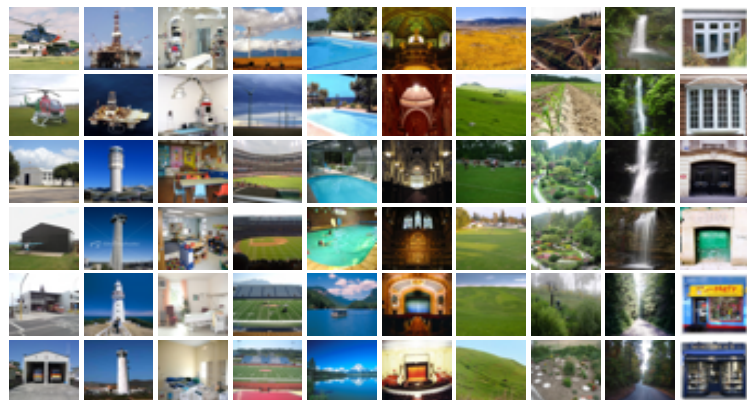
MNIST: Variational versus Gibbs



- Five random initializations from $K=1$, $K=50$, $K=300$ clusters
- Diagonal-covariance Gaussians (change from previous slides)

Clustering Image Patches

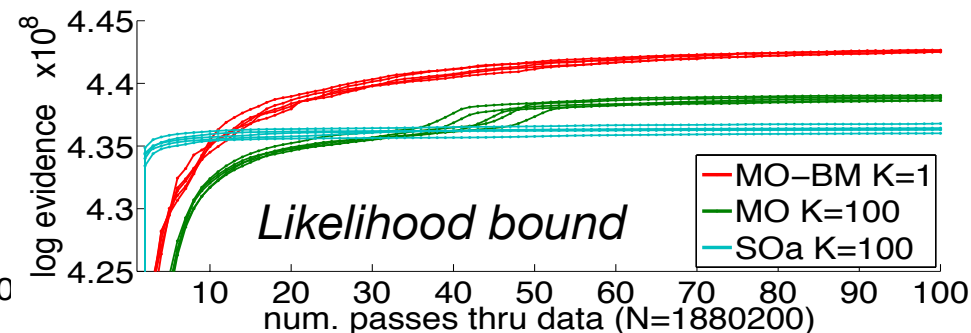
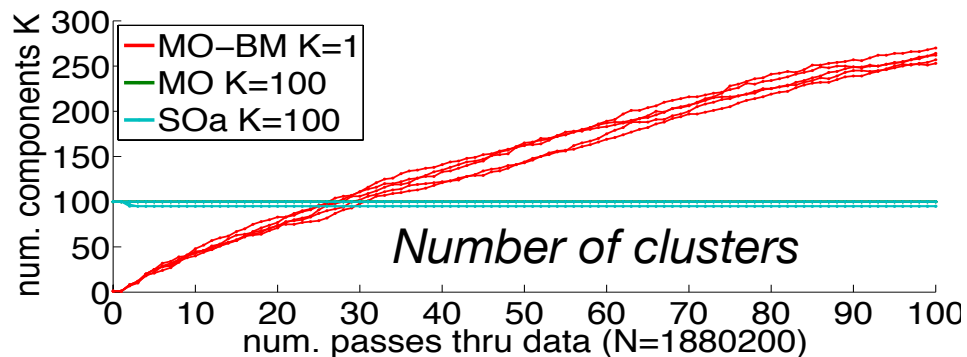
SUN Database of Natural Scene Categories: $N=108,754$



- **Memoized birth-merge** learns a more accurate model with only $K=28$ clusters

8x8 Image Patches (Berkeley Segmentation): $N=1.88$ million

- **Memoized birth-merge** allows growth in model complexity
- Effective performance as density model for image denoising



Memoized Variational Inference for Hierarchical DP Topic Models

Michael Hughes, Dae Il Kim, & E. Sudderth

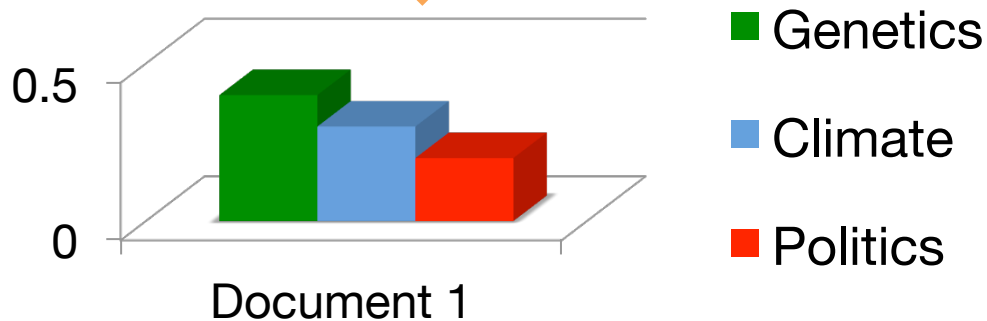
estimation
data density approach em
probability model number set
mixture gaussian posterior bayesian distribution
figure parameters models
log likelihood prior



What are Topic Models?

GOAL: Summarize semantic content of a large document corpus.

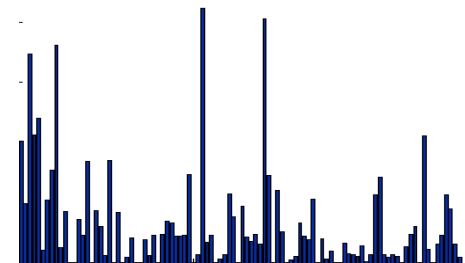
*There are reasons to believe that the **genetics** of an **organism** are likely to shift due to the **extreme changes** in our **climate**. To protect them, our **politicians** must pass **environmental legislation** that can protect our future **species** from becoming **extinct**...*



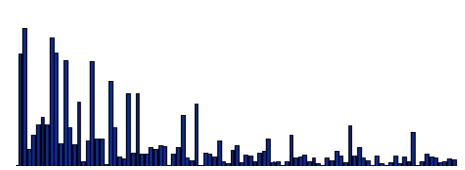
Documents are represented as mixtures of “topics” used with varying frequencies.

Topics are categorical distributions on a (typically large) discrete vocabulary:

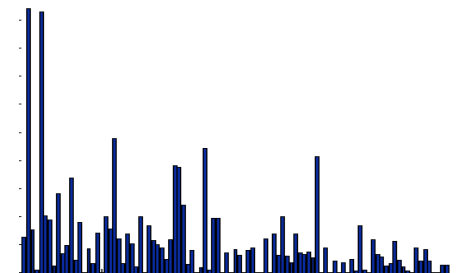
“Genetics”
Topic



“Climate
Change”
Topic



“Politics”
Topic



Hierarchical DP Topic Model

Generalization of Latent Dirichlet Allocation (LDA, Blei 2003) by Teh et al. JMLR 2006.
 Dependent Dirichlet process (DDP, MacEachern 1999) with group-specific weights.

➤ Global topic frequencies and parameters:

$$\beta_k = u_k \prod_{\ell=1}^{k-1} (1 - u_\ell) \quad u_k \sim \text{Beta}(1, \gamma)$$

$$\phi_k \sim \text{Dirichlet}(\lambda_0) \quad (\textit{sparse})$$

➤ For each of D documents (groups):

➤ Topic frequencies: $\pi_d \sim \text{DP}(\alpha\beta)$

Generalized Dirichlet, Connor & Mosimann 1969

$$\pi_{dk} = v_{dk} \prod_{\ell=1}^{k-1} (1 - v_{d\ell})$$

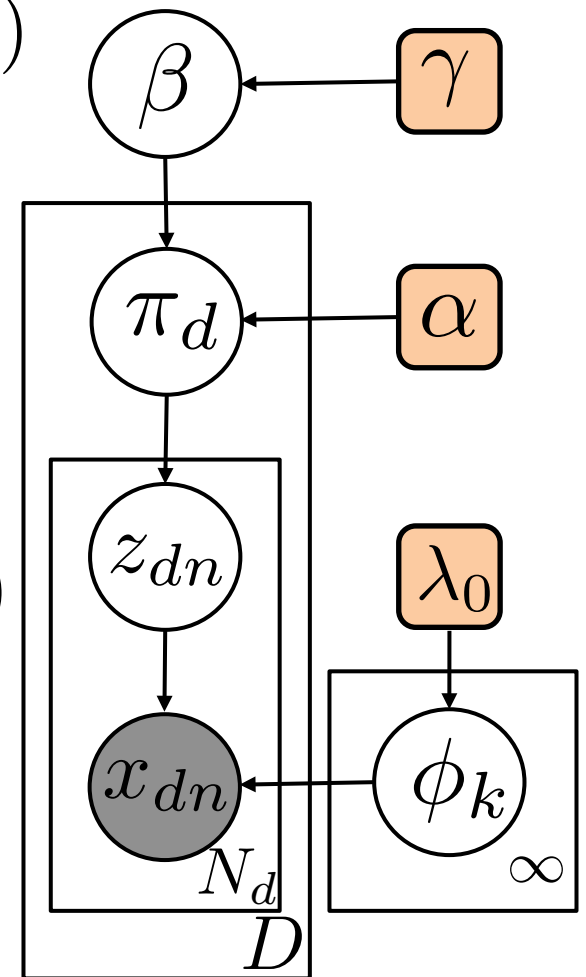
$$v_{dk} \sim \text{Beta}(\alpha_k u_k, \alpha_k (1 - u_k))$$

$$\alpha_k = \alpha \prod_{\ell=1}^{k-1} (1 - v_{d\ell})$$

➤ For each of N_d words in document d :

➤ Topic assignment: $z_{dn} \sim \text{Cat}(\pi_d)$

➤ Observed value: $x_{dn} \sim \text{Cat}(\phi_{z_{dn}})$



Variational Learning of HDP Topics

$$q(z_{dn}) = \text{Cat}(z_{dn} \mid r_{dn1}, r_{dn2}, \dots, r_{dnK}, 0, 0, 0, \dots) \quad \text{for some } K > 0$$

Update Document Distributions: For $k \leq K$,

$$r_{dnk} \propto \exp(\mathbb{E}_q[\log \pi_{dk}(v_d)] + \mathbb{E}_q[\log p(x_{dn} \mid \phi_k)])$$

$$\mathbb{E}_q[\log \pi_{dk}(v_d)] = \mathbb{E}_q[\log(v_{dk})] + \sum_{\ell=1}^{k-1} \mathbb{E}_q[\log(1 - v_{d\ell})]$$

- Closed form update for beta stick-breaking weights
- Local iteration between assignments and weights

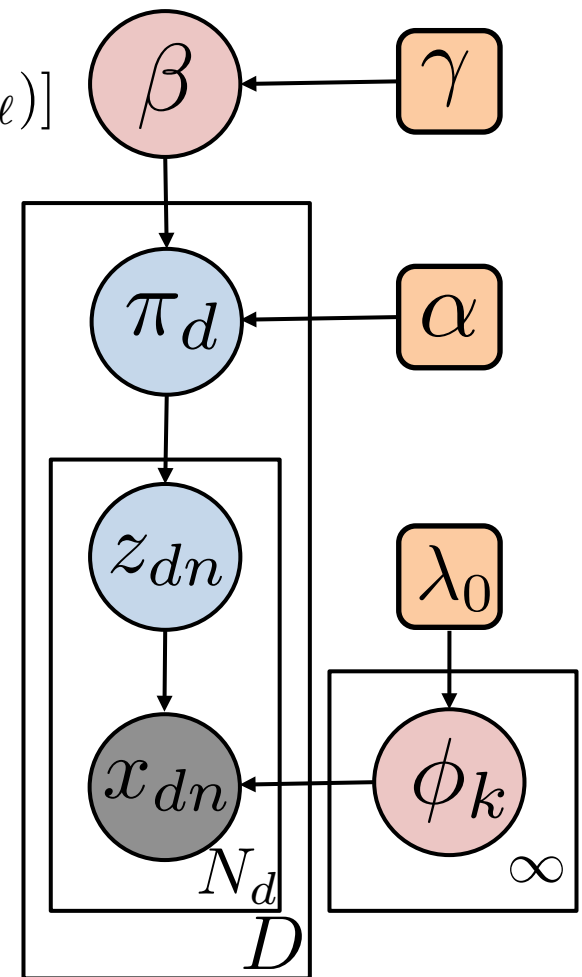
Update Global Parameters:

$$q(\phi_k) = \text{Dir}(\phi_k \mid \lambda_0 + s_k^0)$$

$$s_k^0 \leftarrow \sum_{d=1}^D \sum_{n=1}^{N_d} r_{dnk} t(x_{dn})$$

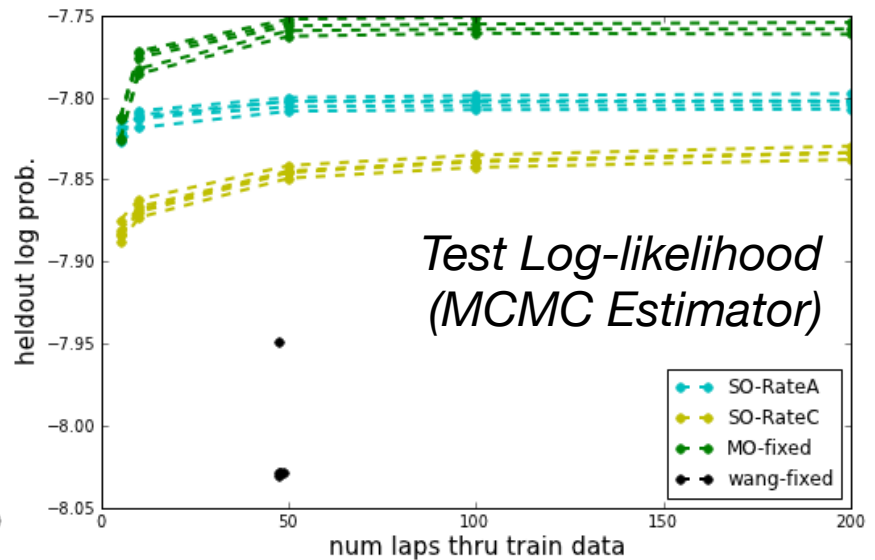
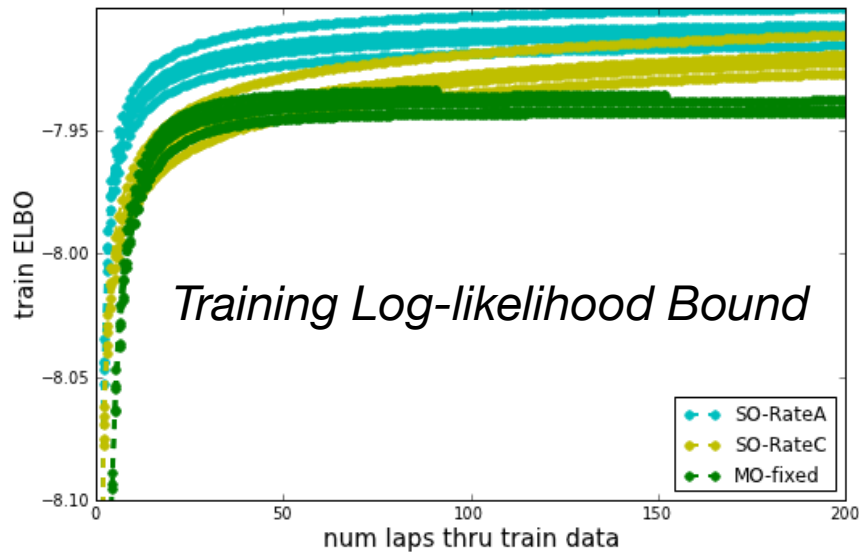
- Closed form for topic-specific word distributions
- **Beta normalization constants have non-conjugate dependence on topic frequencies, requires additional bound and numerical optimization**

Iterate: Batch, Stochastic, or Memoized

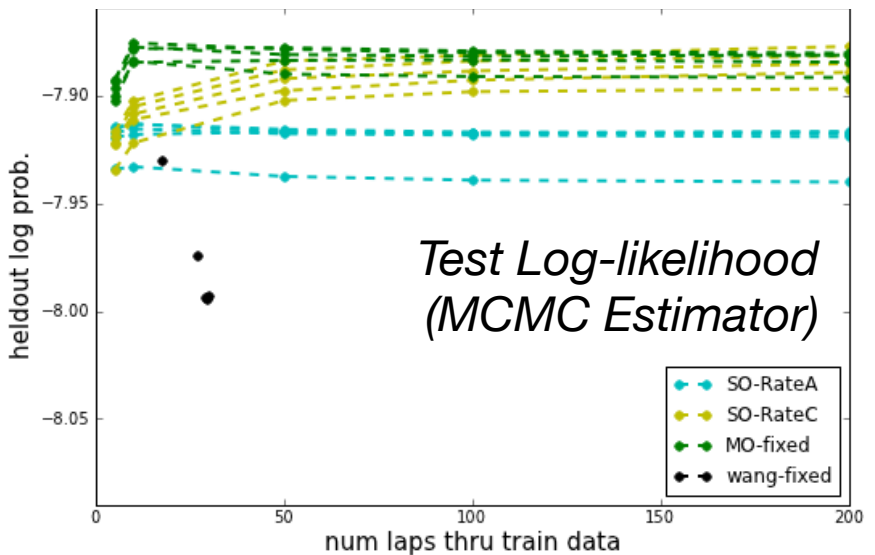
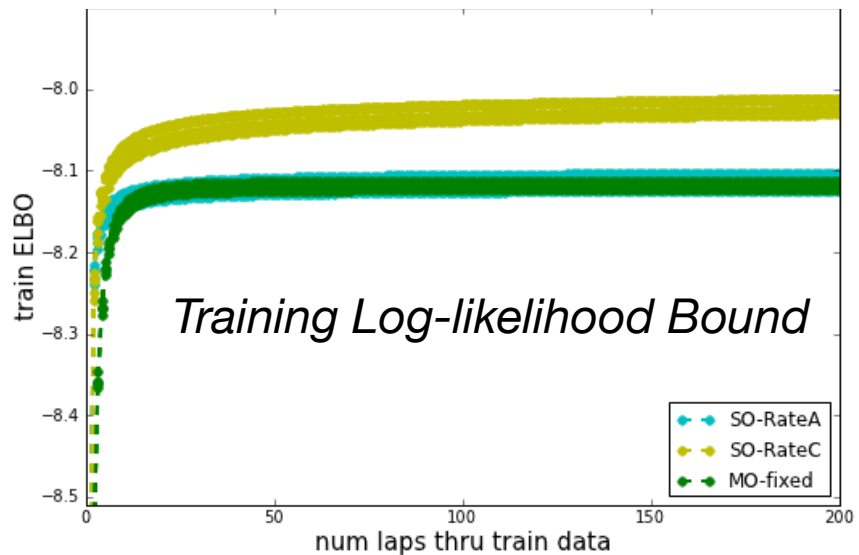


Analysis of Document Corpora

NIPS Conference
(D=1392)



Huffington Post
(D=3271)

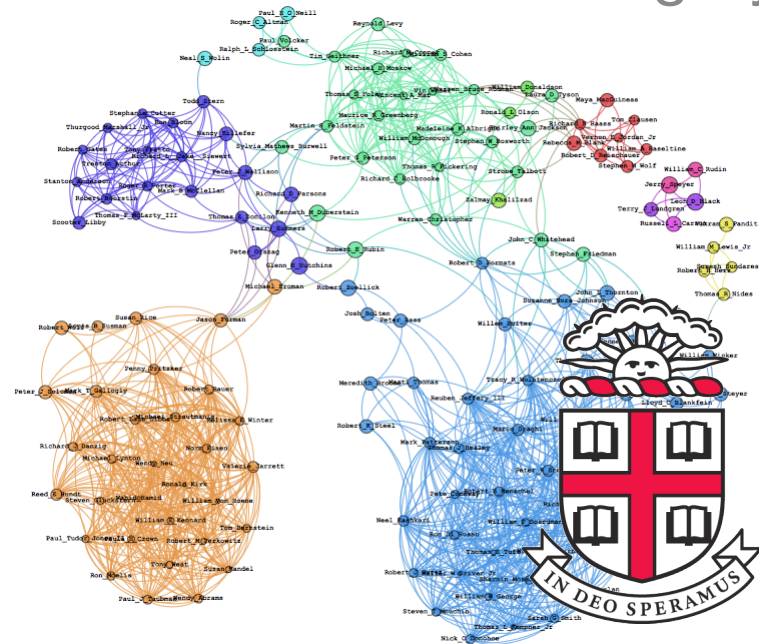


- Variational: Memoized versus stochastic rate A, stochastic rate C
- Baseline: Stochastic variational on “expanded” HDP (Wang et al. 2011)

Stochastic Variational Inference for Hierarchical DP Relational Models

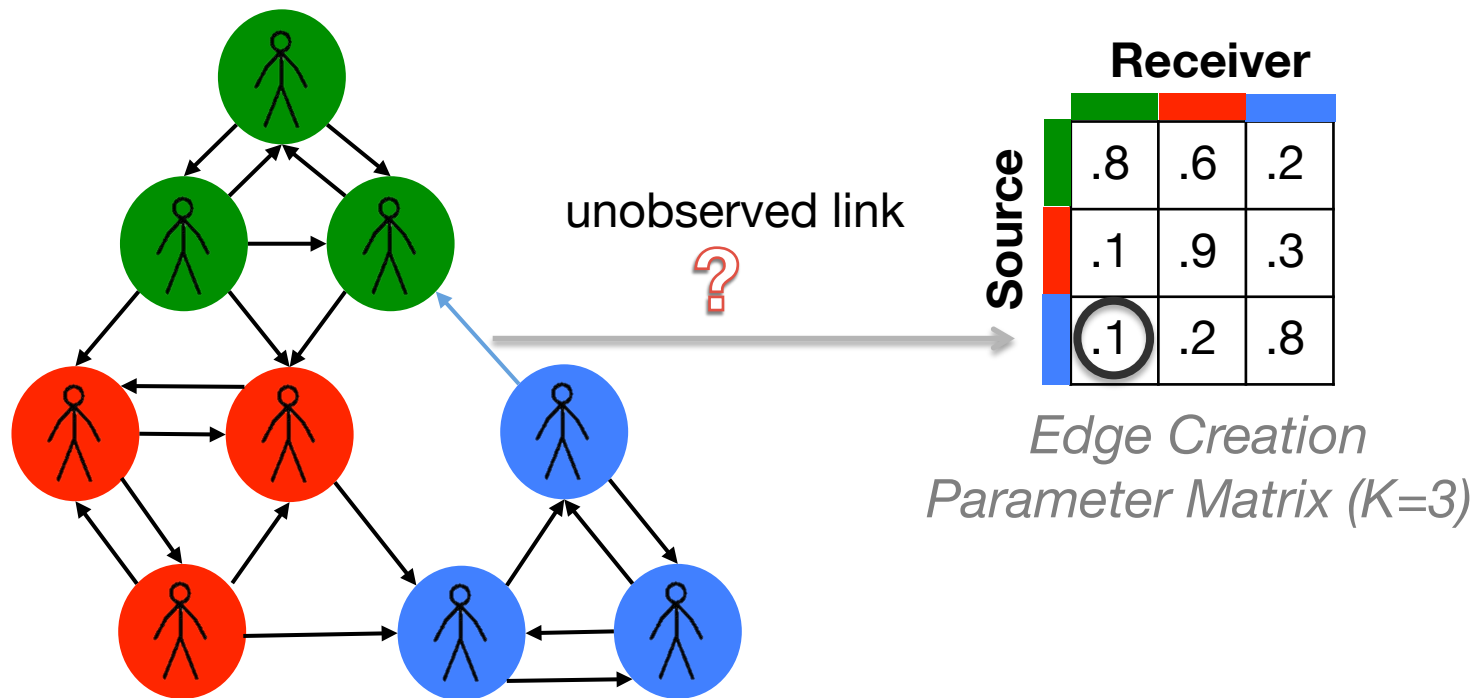
D. Kim, P. Gopalan, D. Blei, & E. Sudderth

2013 Conference on Neural Information Processing Systems



What are Relational Models?

GOAL: *Unsupervised community discovery from observed relationships.*

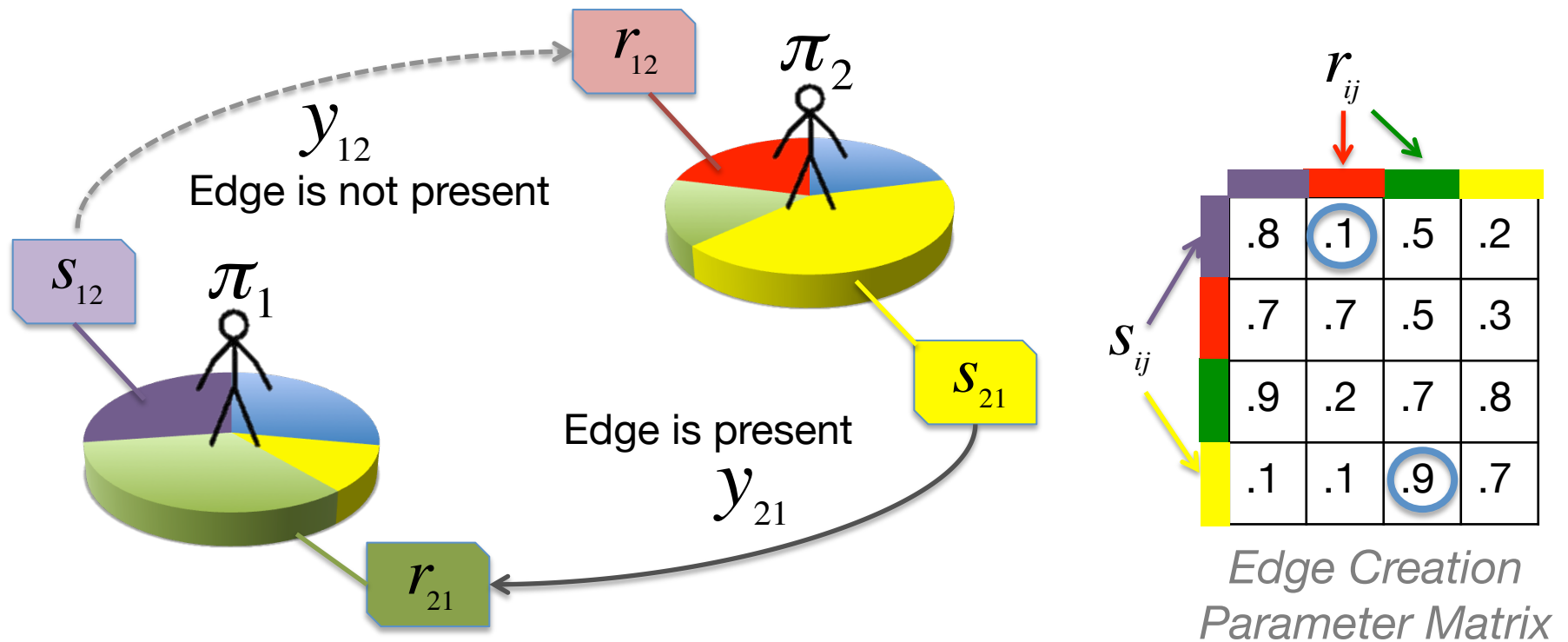


Stochastic Block Model: (Wang et al., JASA 1987)

- Assign each node to *one* latent block/community
- Predict edge presence or absence from block assignments of *source* and *receiver* nodes

Mixed Membership Blockmodels

Parametric mixed membership stochastic blockmodel, Airoldi et al. JMLR 2008



- Source Community Assignment
- Receiver Community Assignment
- Community Link Probability
- Binary Edge Indicator

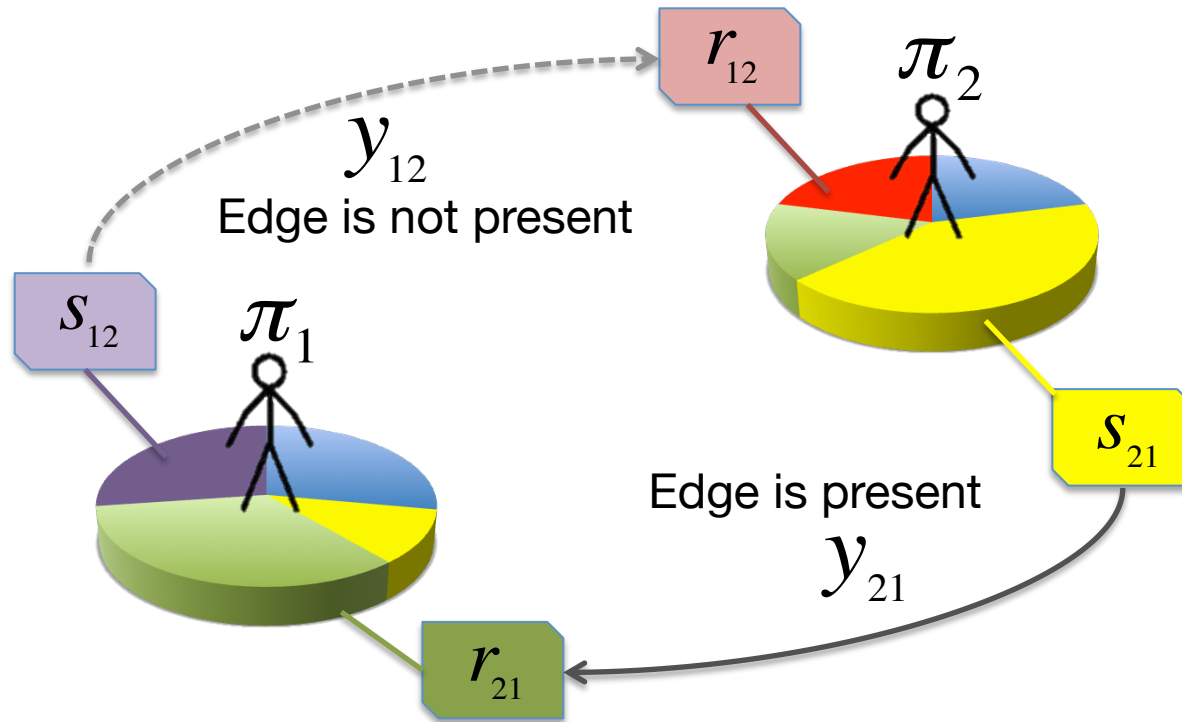
$$s_{ij} \sim \text{Cat}(\pi_i)$$

$$r_{ij} \sim \text{Cat}(\pi_j)$$

$$\phi_{kl} \sim \text{Beta}(\tau_a, \tau_b)$$

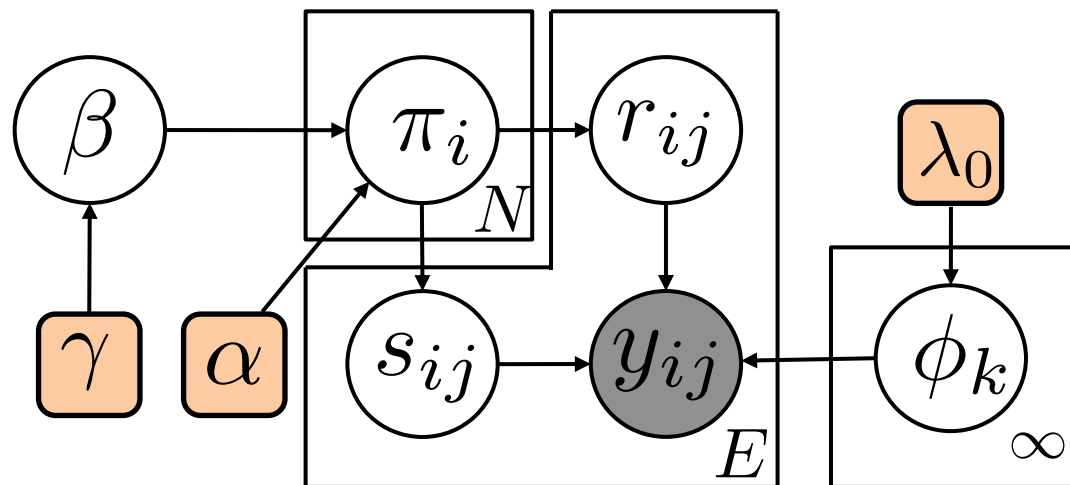
$$y_{ij} \sim \text{Bern}(s_{ij} \phi r_{ij}^T)$$

HDP Relational Models



	.8	.1	.5	.2
.7	.7	.5	.3	
.9	.2	.7	.8	
.1	.1	.9	.7	

Edge Creation
Parameter Matrix



$$s_{ij} \sim \text{Cat}(\pi_i)$$

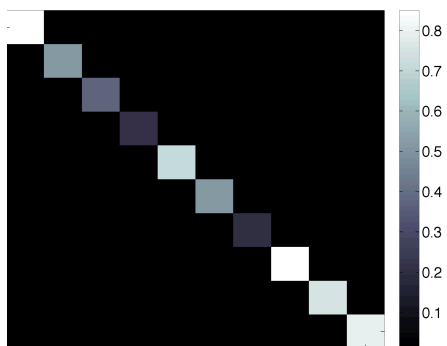
$$r_{ij} \sim \text{Cat}(\pi_j)$$

$$\phi_{kl} \sim \text{Beta}(\tau_a, \tau_b)$$

$$y_{ij} \sim \text{Bern}(s_{ij} \phi_{r_{ij}}^T)$$

Variational Learning of Relations

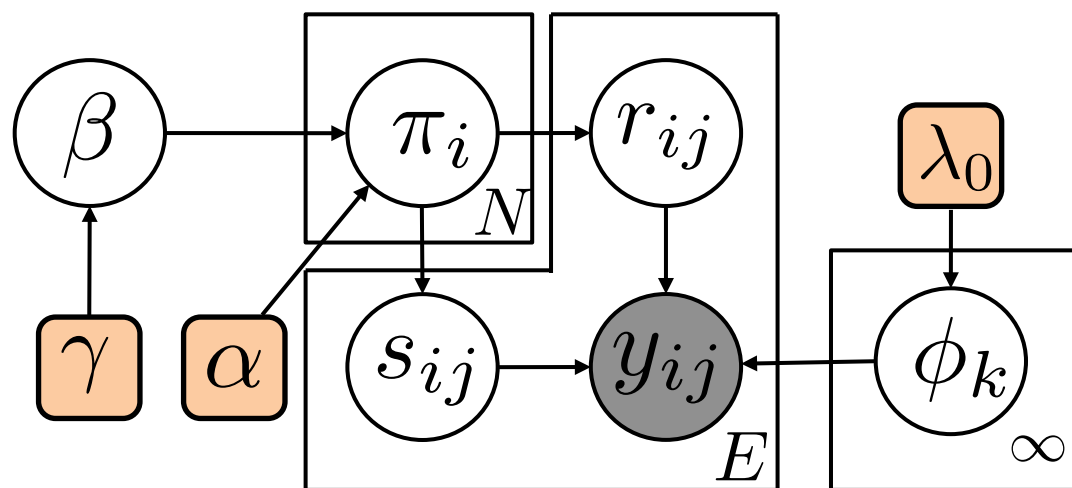
Assortative Likelihoods:



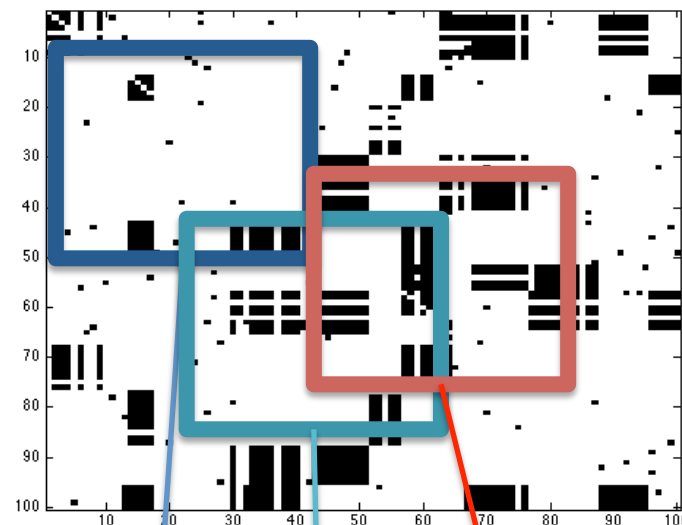
$$p(y_{ij} = 1 \mid s_{ij} = r_{ij} = k) = \phi_k$$

$$p(y_{ij} = 1 \mid s_{ij} \neq r_{ij}) = \epsilon$$

➔ $\mathcal{O}(K)$ storage & computation for distribution on K^2 community pairs



Stochastic Variational:

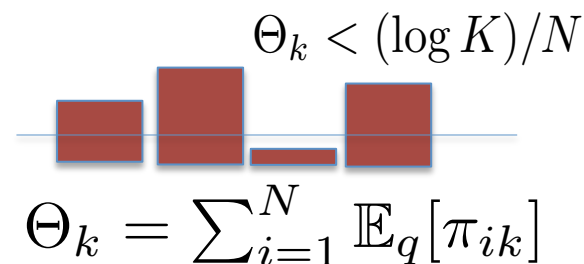


Mini-Batch #1

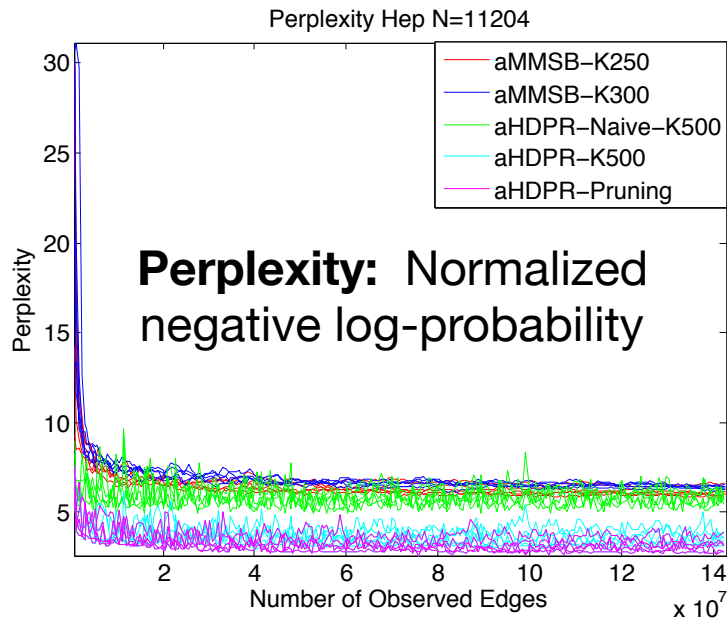
Mini-Batch #2

Mini-Batch #3

Variational Pruning:



Analysis of Collaboration Networks

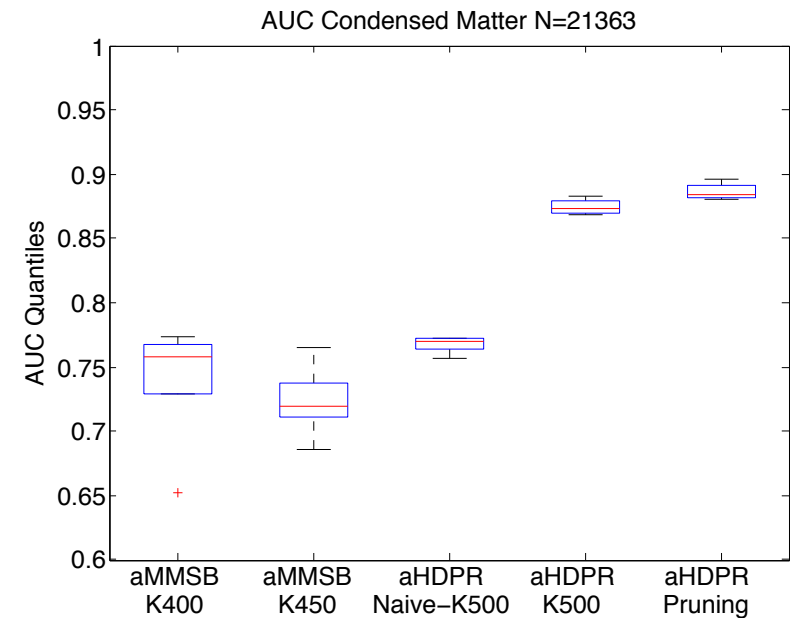
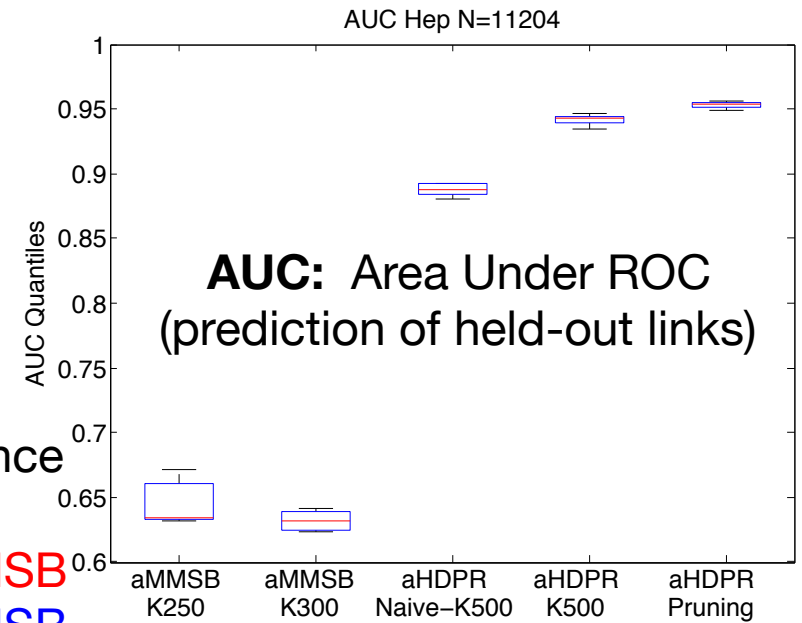
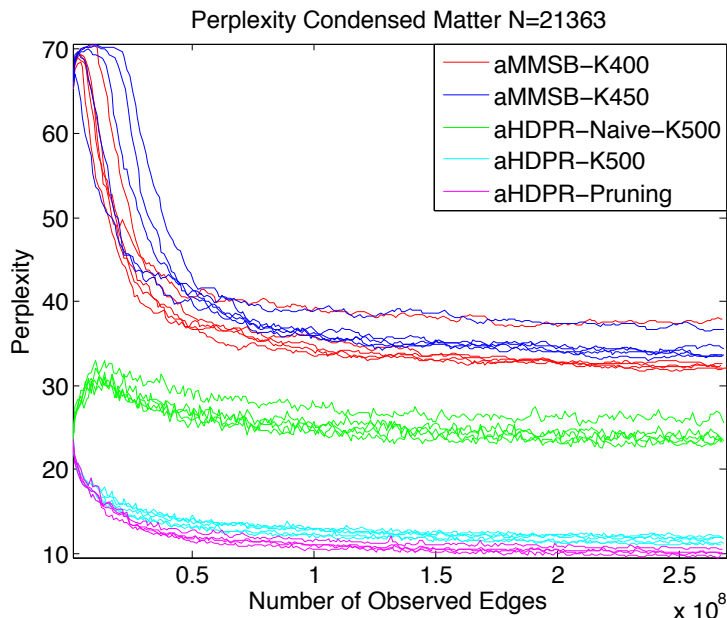


High Energy Physics (HEP)
N=11,204

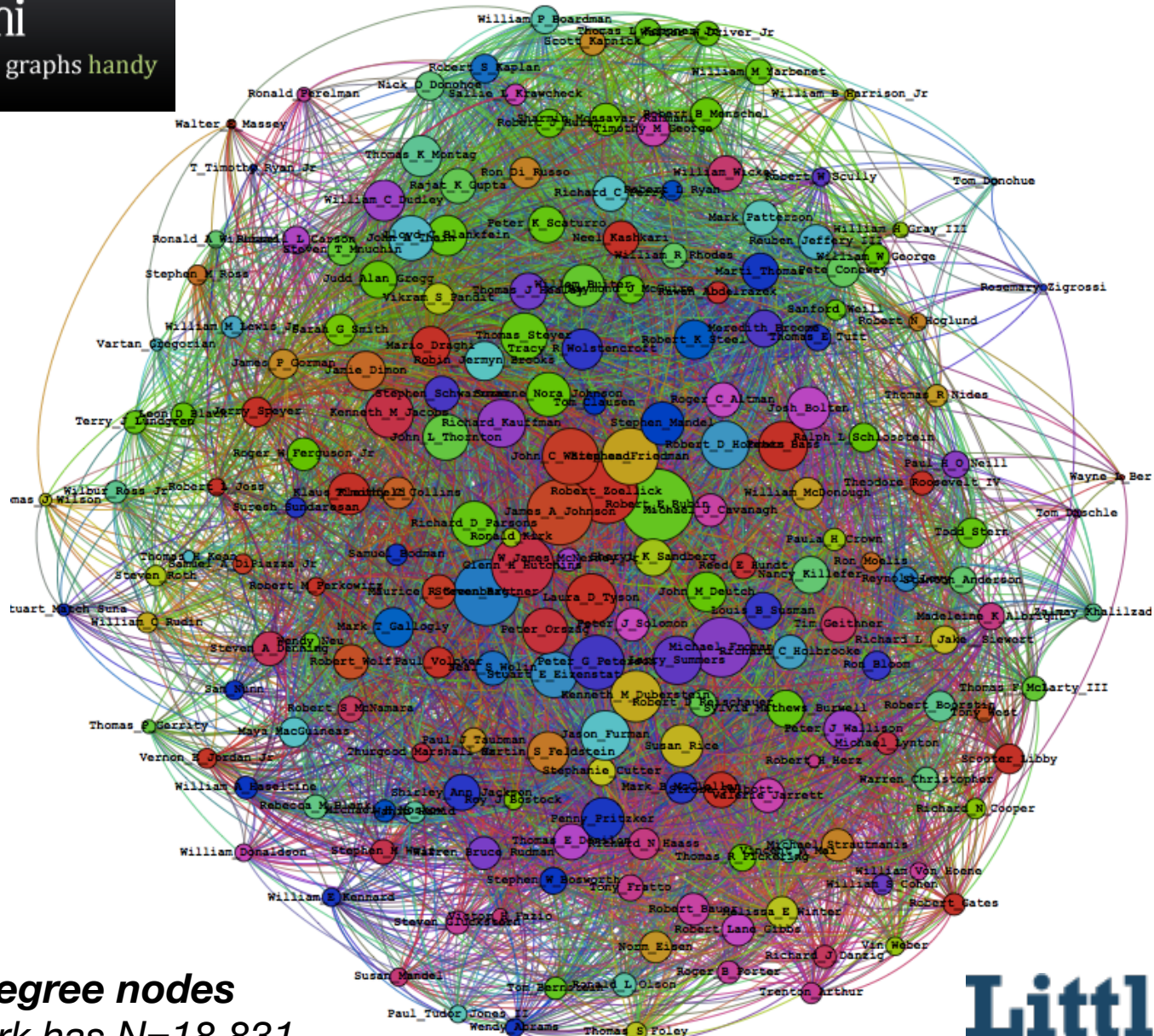
Stochastic inference & model variants:

- Parametric MMSB
- Parametric MMSB
- HDP naïve
- HDP blocked
- HDP pruned

Condensed Matter Physics
N=21,363



LittleSis Network: Raw Data



Top 200 degree nodes
Full network has N=18,831

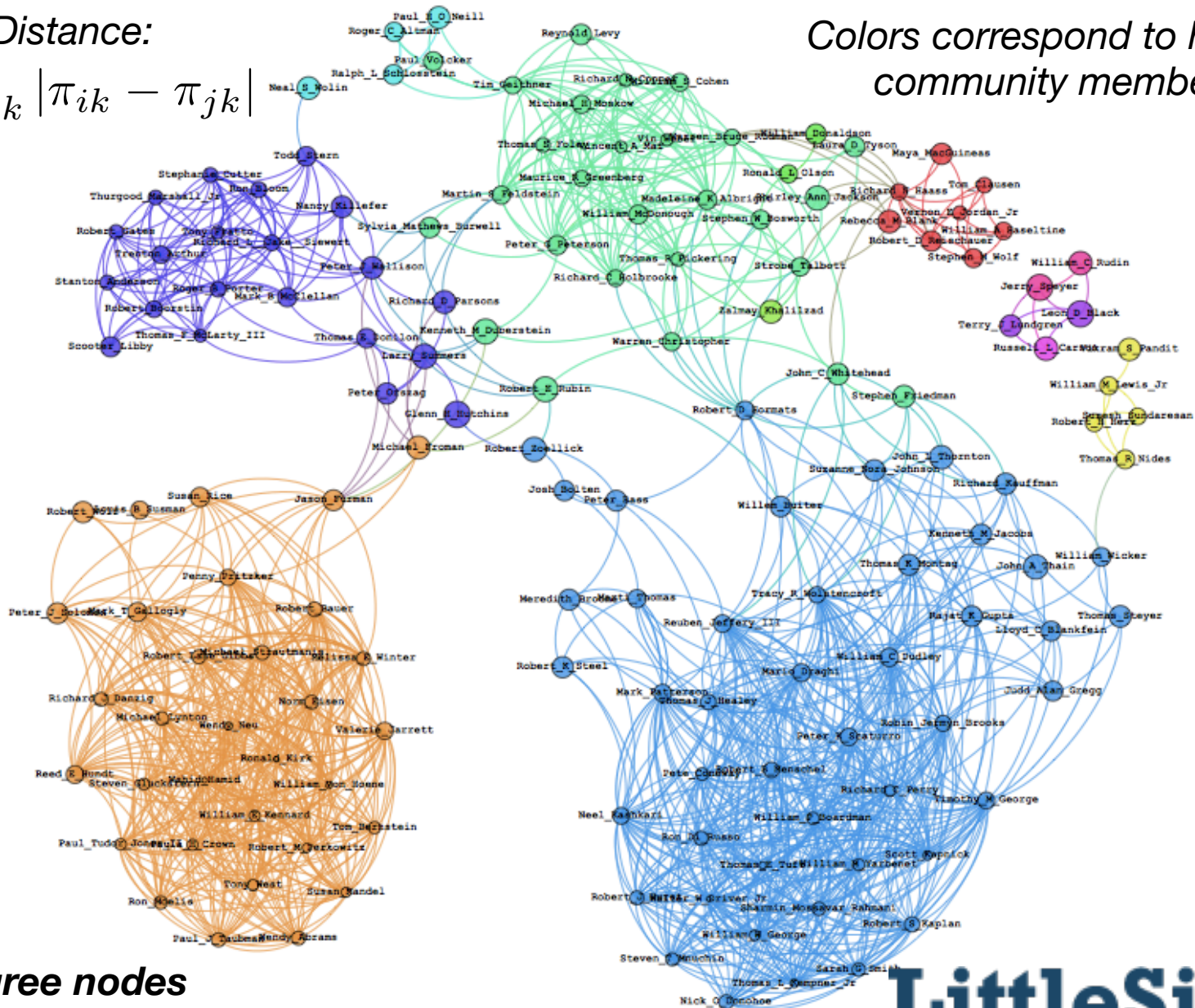


LittleSis Network Communities

Community Distance:

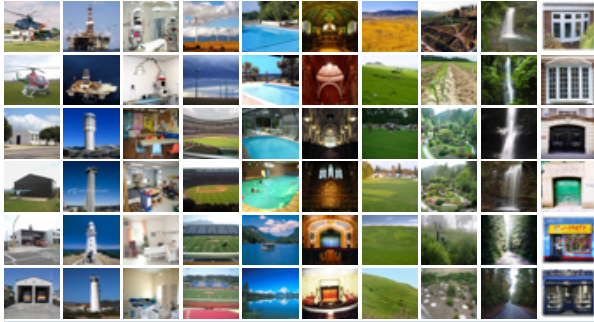
$$D_{ij} = \frac{1}{2} \sum_k |\pi_{ik} - \pi_{jk}|$$

Colors correspond to highest community memberships

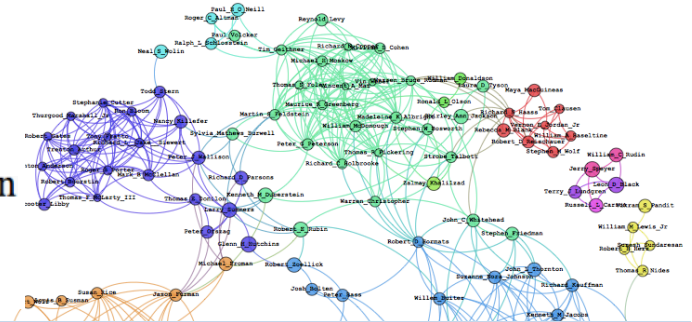


Top 200 degree nodes
Full network has N=18,831



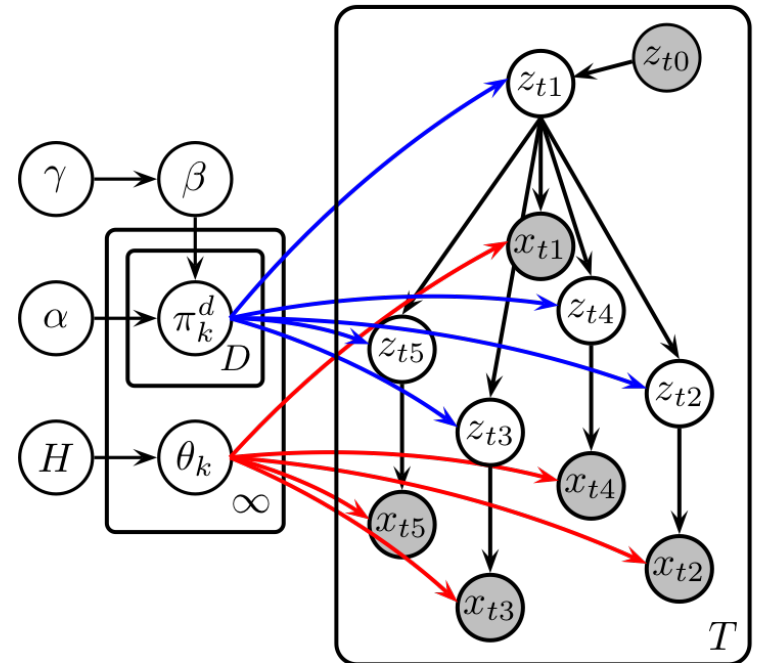


estimation data density approach em probability model number set mixture gaussian posterior bayesian distribution figure parameters models log likelihood prior



Reliable Variational Learning for Hierarchical Dirichlet Processes

- **Scalable:** Large-scale learning via stochastic or memoized updates
- **Reliable:** Birth-merge recovers structure informed by model & data, not inference algorithm limitations
- **Flexible:** Designed to be broadly applicable: space, time, scale, ...



BNPy: Bayesian Nonparametric Learning in Python

Erik Sudderth @ Brown CS:

<http://cs.brown.edu/~sudderth/>