

The Nested Chinese Restaurant Process and Bayesian Nonparametric Inference of Topic Hierarchies

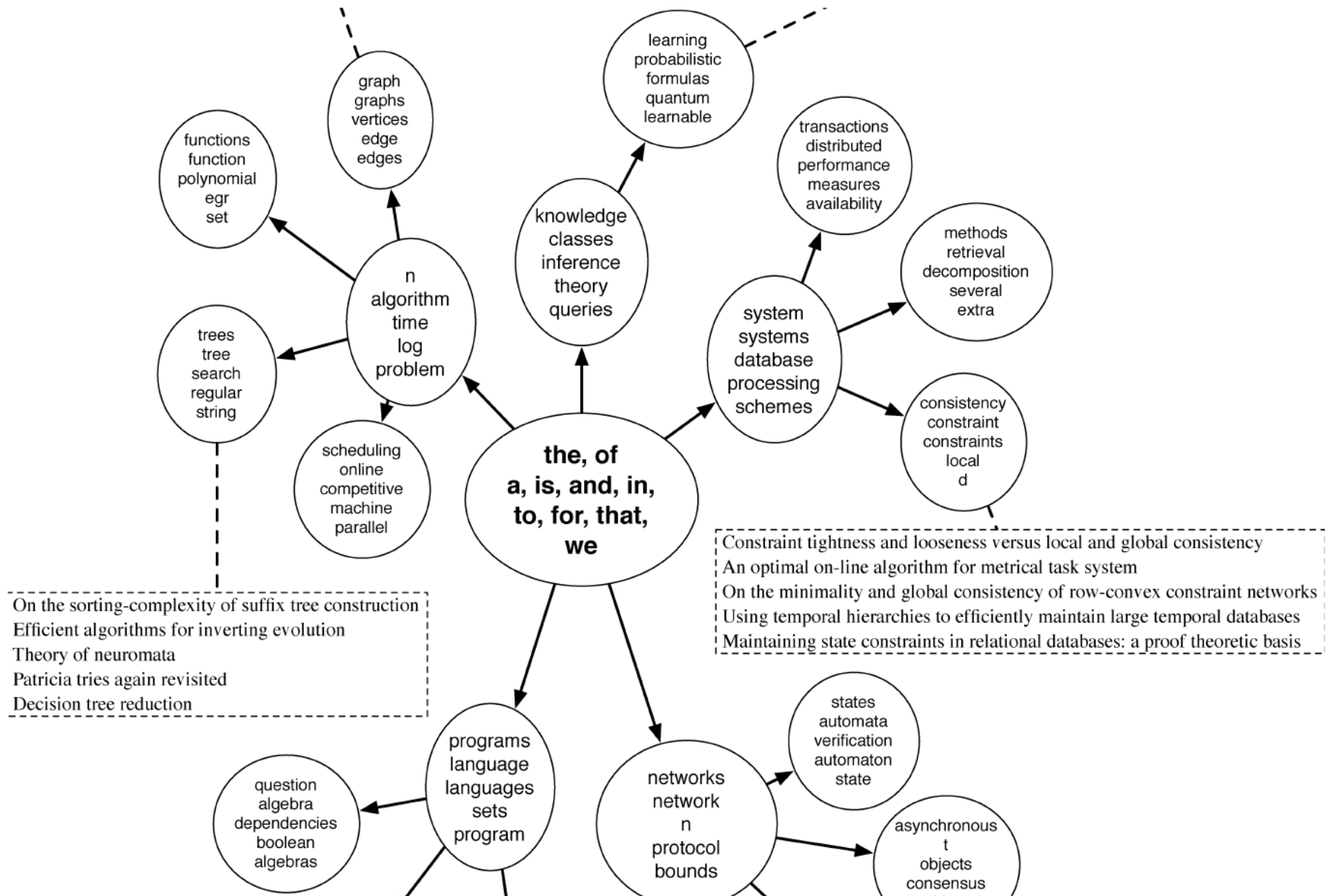
David M. Blei, Thomas L. Griffiths, and Michael I. Jordan
JACM, January 2010

Michael Bryant
Jixiong Wang

Motivation

Problem:

- To learn topic models for collections of text, images and other semi-structured corpora.
- The original topic models treat topics as a “flat” set of probability distributions, with no direct **relationship** between one topic and another. They fail to indicate the level of **abstraction** of a topic, or how the various topics are related.
- We want an algorithm to both find useful sets of topics and learn to organize the topics according to a hierarchy in which **more abstract topics are near the root of the hierarchy and more concrete topics are near the leaves.**

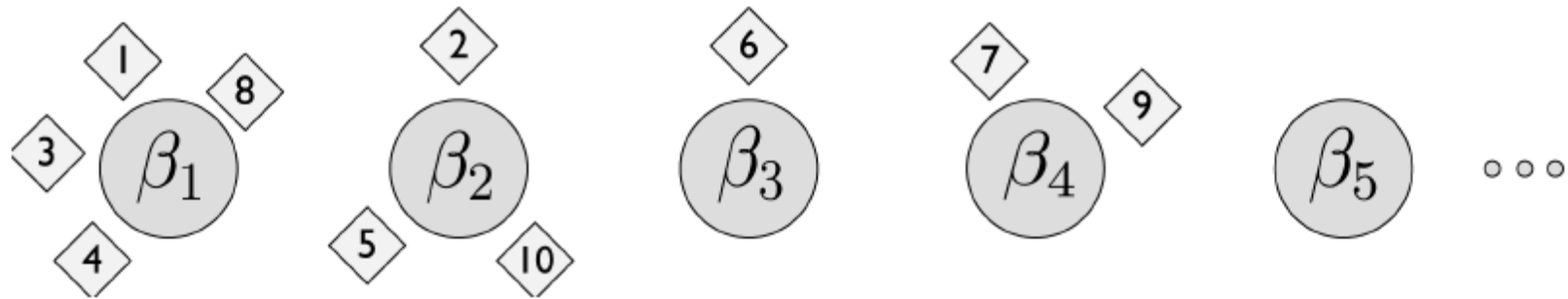


Motivation

Problem:

- To learn topic models for collections of text, images and other semi-structured corpora.
- The original topic models treat topics as a “flat” set of probability distributions, with no direct **relationship** between one topic and another. They fail to indicate the level of **abstraction** of a topic, or how the various topics are related.
- We want an algorithm to both find useful sets of topics and learn to organize the topics according to a hierarchy in which **more abstract topics are near the root of the hierarchy and more concrete topics are near the leaves**.
- While a classical unsupervised analysis might require the topology of the hierarchy to be chosen in advance, we want the approach to place high probability on those hierarchies that best explain the data. We need a **distribution on topologies**.
- We wish to allow this distribution to have its support on arbitrary topologies, -- there should be no limitations such as a maximum depth or maximum branching factor. **BNP is needed**.

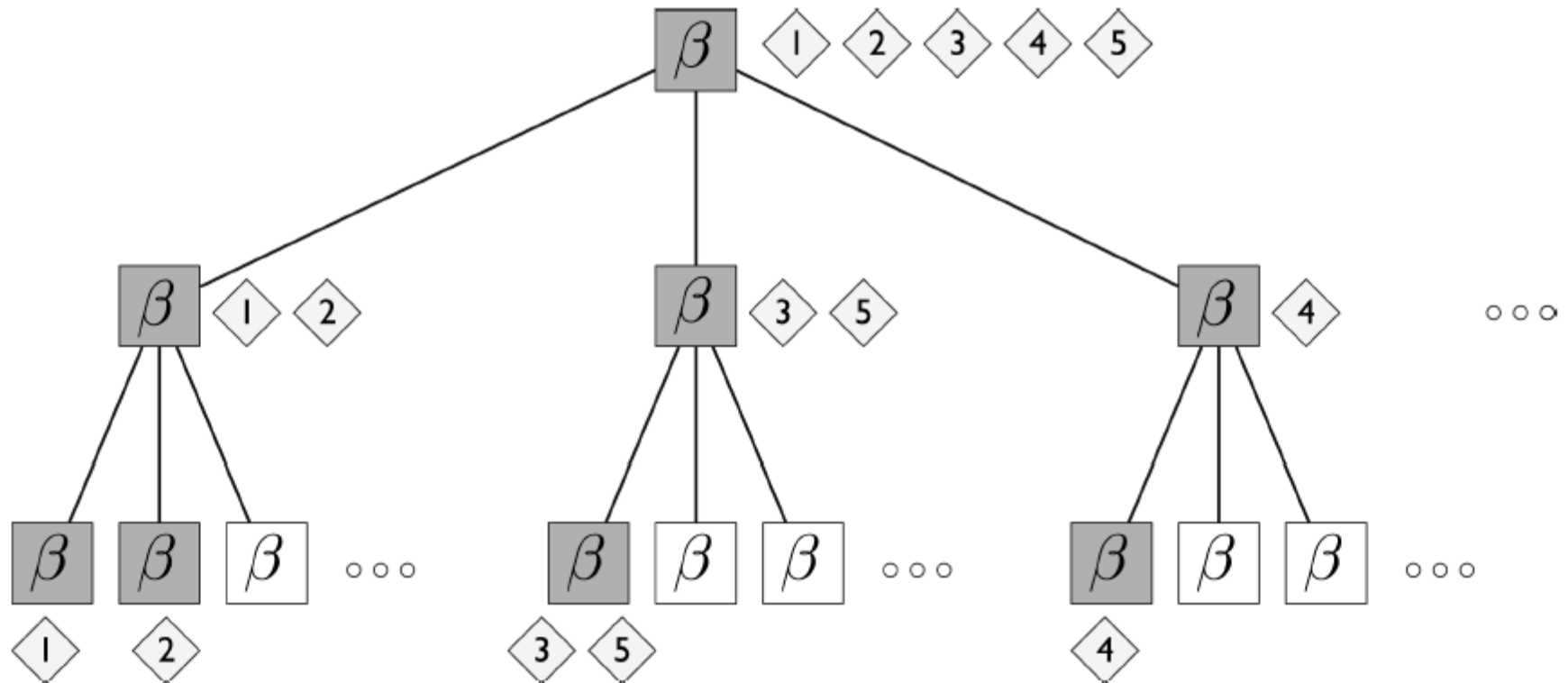
Chinese Restaurant Process



$$p(\text{occupied table } i \mid \text{previous customers}) = \frac{n_i}{\gamma + n - 1}$$

$$p(\text{next unoccupied table} \mid \text{previous customers}) = \frac{\gamma}{\gamma + n - 1}$$

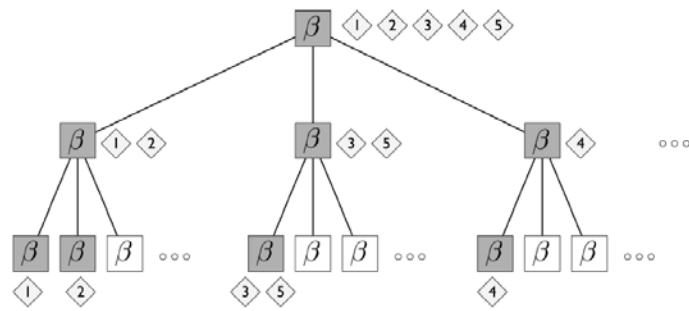
Nested Chinese Restaurant Process



$$p(\text{occupied table } i \mid \text{previous customers}) = \frac{n_i}{\gamma + n - 1}$$

$$p(\text{next unoccupied table} \mid \text{previous customers}) = \frac{\gamma}{\gamma + n - 1}$$

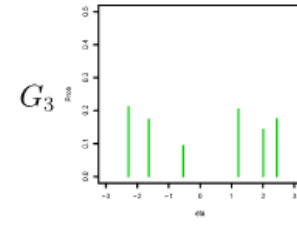
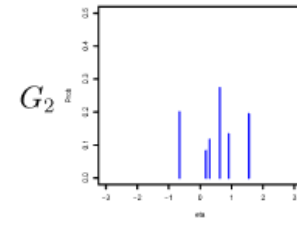
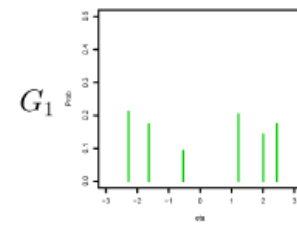
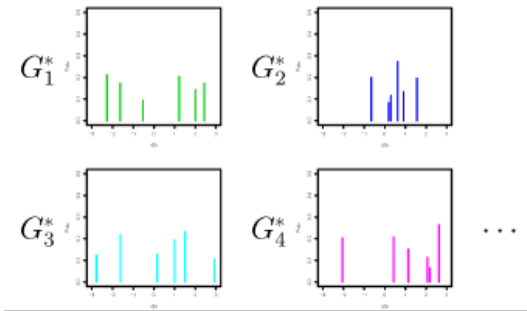
nCRP vs nDP



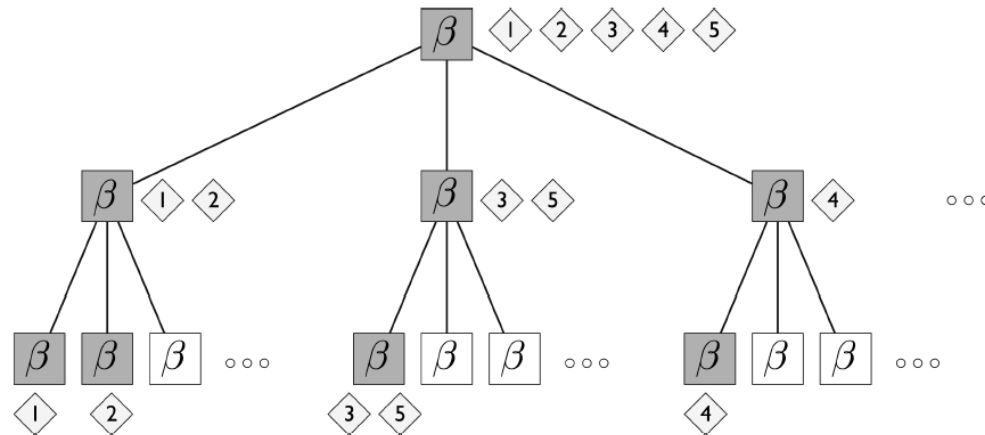
NDP

$$G_j \sim Q$$

$$Q \sim \text{DP}(\alpha \text{DP}(\beta H))$$



Hierarchical Latent Dirichlet Allocation



for each table $k \in \mathcal{T}$ in the infinite tree do

- Draw a topic $\beta_k \sim \text{Dirichlet}(\eta)$

end for

for each document $d \in \{1, 2, \dots, D\}$ do

- Draw $c_d \sim \text{nCRP}(\gamma)$

- Draw a distribution over levels in the tree, $\theta_d | \{m, \pi\} \sim \text{GEM}(m, \pi)$

for each word in document d do

- Choose level $Z_{d,n} | \theta_d \sim \text{Discrete}(\theta_d)$

- Choose word $W_{d,n} | \{z_{d,n}, c_d, \beta\} \sim \text{Discrete}(\beta_{c_d}[z_{d,n}])$

end for

end for

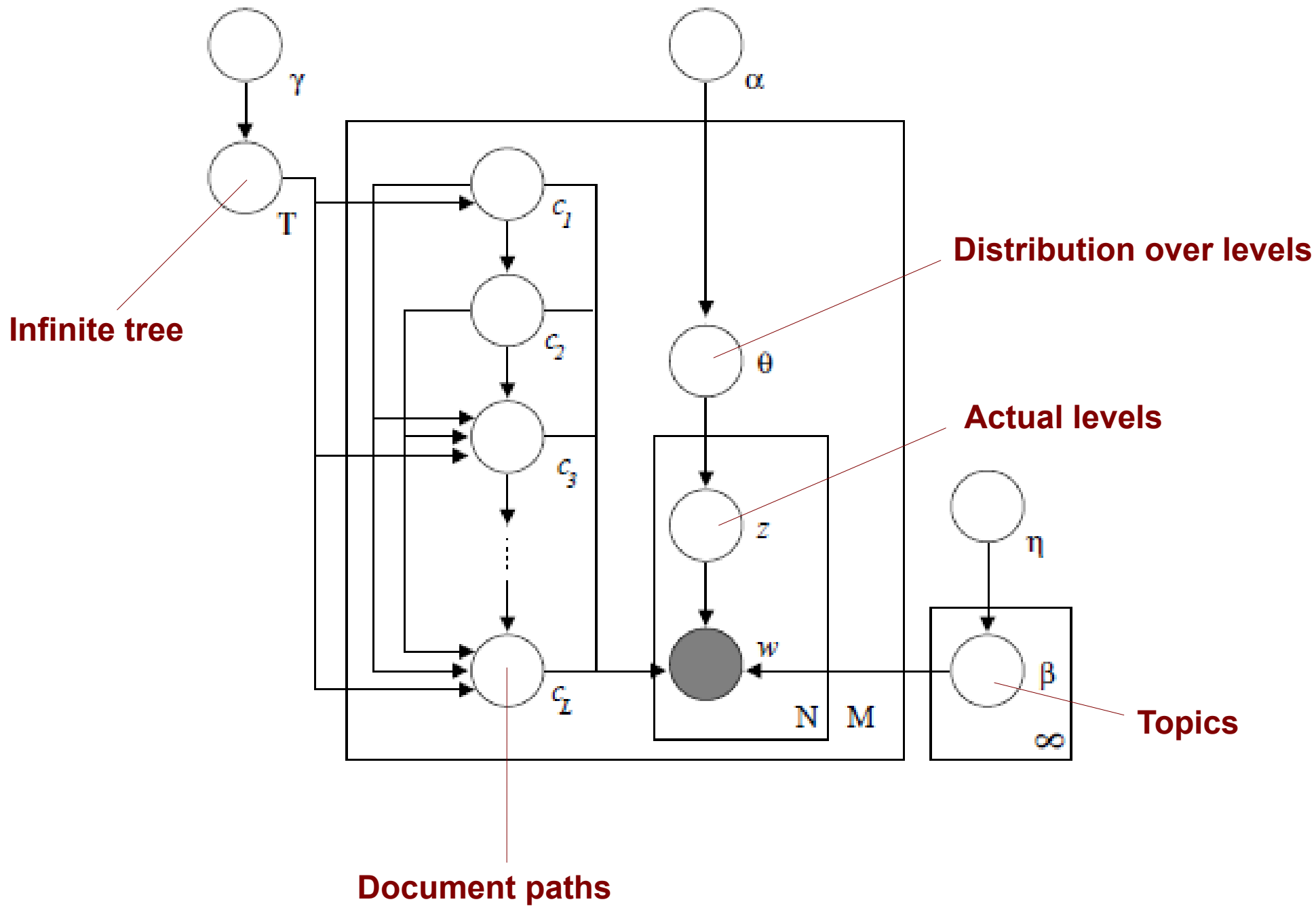
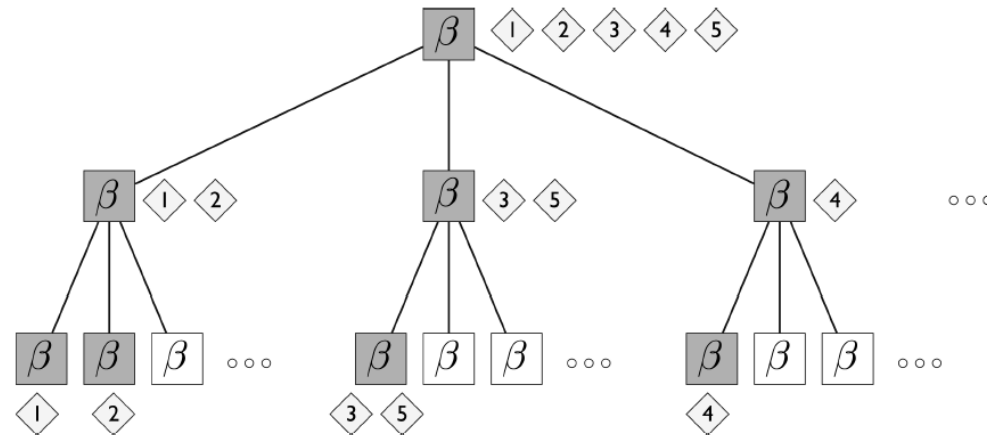


Figure from Blei, et al 2003

Hierarchical Latent Dirichlet Allocation



for each table $k \in \mathcal{T}$ in the infinite tree do

- Draw a topic $\beta_k \sim \text{Dirichlet}(\eta)$

end for

for each document $d \in \{1, 2, \dots, D\}$ do

- Draw $c_d \sim \text{nCRP}(\gamma)$

- Draw a distribution over levels in the tree, $\theta_d | \{m, \pi\} \sim \text{GEM}(m, \pi)$

for each word in document d do

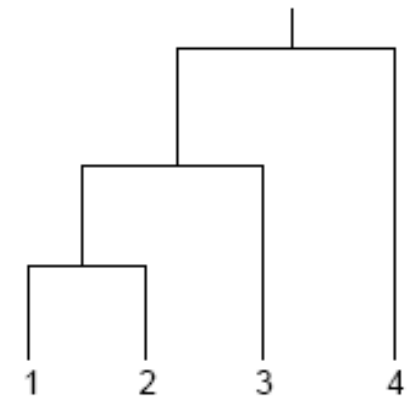
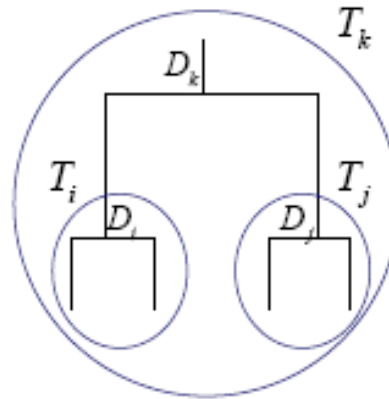
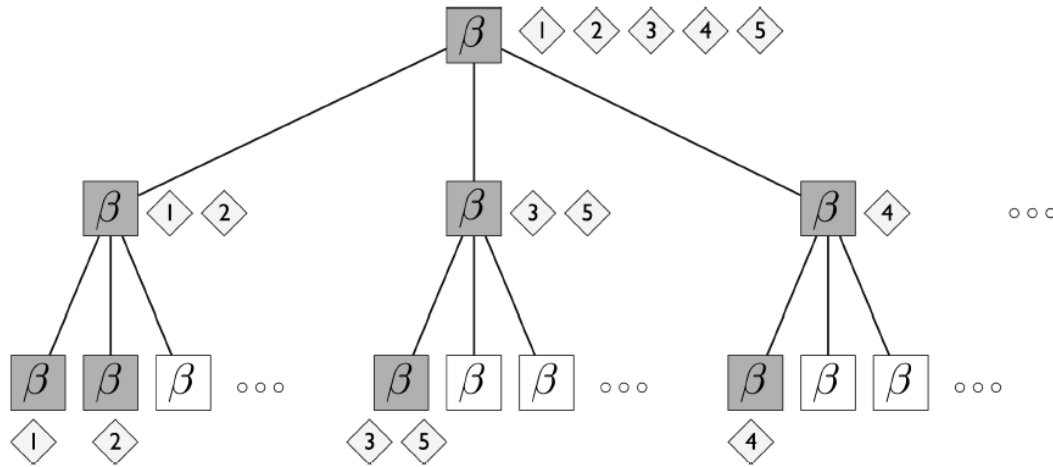
- Choose level $Z_{d,n} | \theta_d \sim \text{Discrete}(\theta_d)$

- Choose word $W_{d,n} | \{z_{d,n}, c_d, \beta\} \sim \text{Discrete}(\beta_{c_d}[z_{d,n}])$

end for

end for

hLDA vs Hierarchical Clustering



hLDA vs HDP-LDA

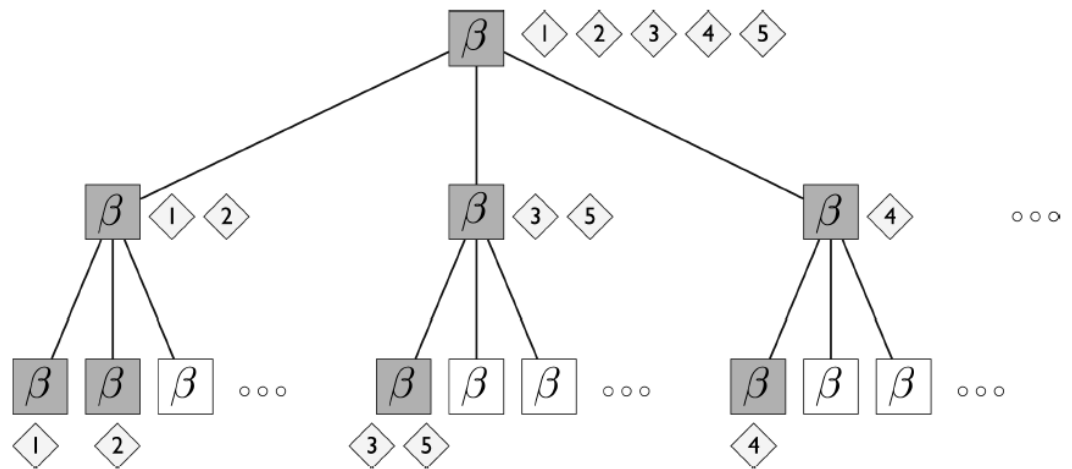
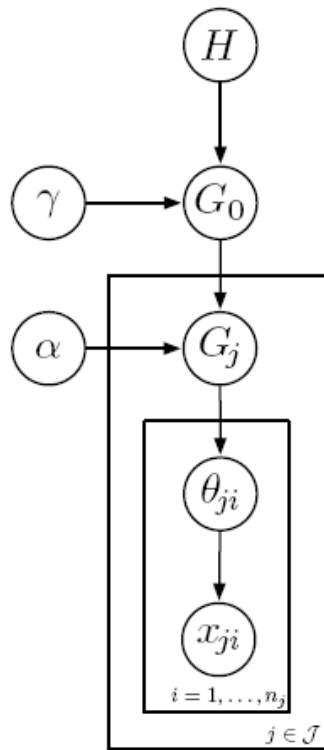
$$G_0 | \gamma, H \sim \text{DP}(\gamma, H)$$

$$G_j | \alpha, G_0 \sim \text{DP}(\alpha, G_0)$$

$$\theta_{ji} | G_j \sim G_j$$

$$x_{ji} | \theta_{ji} \sim F_{\theta_{ji}},$$

for each document $j \in \mathcal{J}$
 for each word $i = 1, \dots, n_j$



Notation

$\mathbf{C}_{1:D}$ = Paths for documents $1, \dots, D$

$\mathbf{Z}_{1:D}$ = Level assignments for documents $1, \dots, D$

$\mathbf{Z}_{-(d,n)}$ = All level assignments excluding that for word n in doc d

$\mathbf{Z}_{d,-n}$ = All level assignments for doc d excluding word n

γ = nCRP hyperparameter

η = Topic hyperparameter

m, π = Level distribution hyperparameters

$\mathbf{W}_{1:D}$ = Words for documents $1, \dots, D$

Inference for hLDA

- Given the model and a corpus, we want to estimate the distribution of (some of) the latent variables. In math words:

$$p(\mathbf{c}_{1:D}, \mathbf{z}_{1:D} | \gamma, \eta, m, \pi, \mathbf{w}_{1:D})$$

- Approximate this posterior by *Collapsed Gibbs Sampling*.
- Iterate between:

(1) Sampling level assignments:

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{c}, \mathbf{w}, m, \pi, \eta) \propto p(z_{d,n} | \mathbf{z}_{d,-n}, m, \pi) p(w_{d,n} | \mathbf{z}, \mathbf{c}, \mathbf{w}_{-(d,n)}, \eta)$$

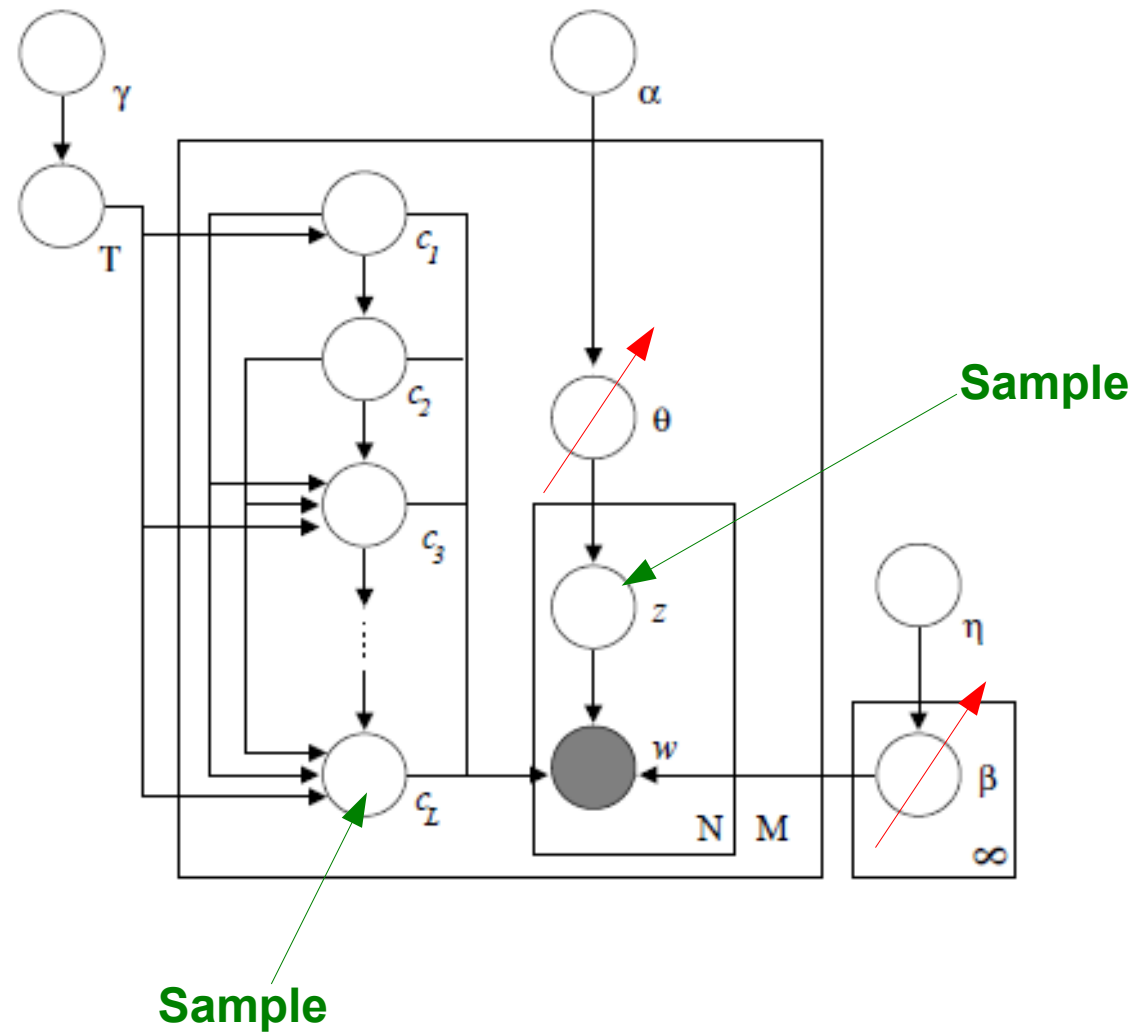
(2) Sampling paths:

$$p(\mathbf{c}_d | \mathbf{w}, \mathbf{c}_{-d}, \mathbf{z}, \eta, \gamma) \propto p(\mathbf{c}_d | \mathbf{c}_{-d}, \gamma) p(\mathbf{w}_d | \mathbf{c}, \mathbf{w}_{-d}, \mathbf{z}, \eta)$$

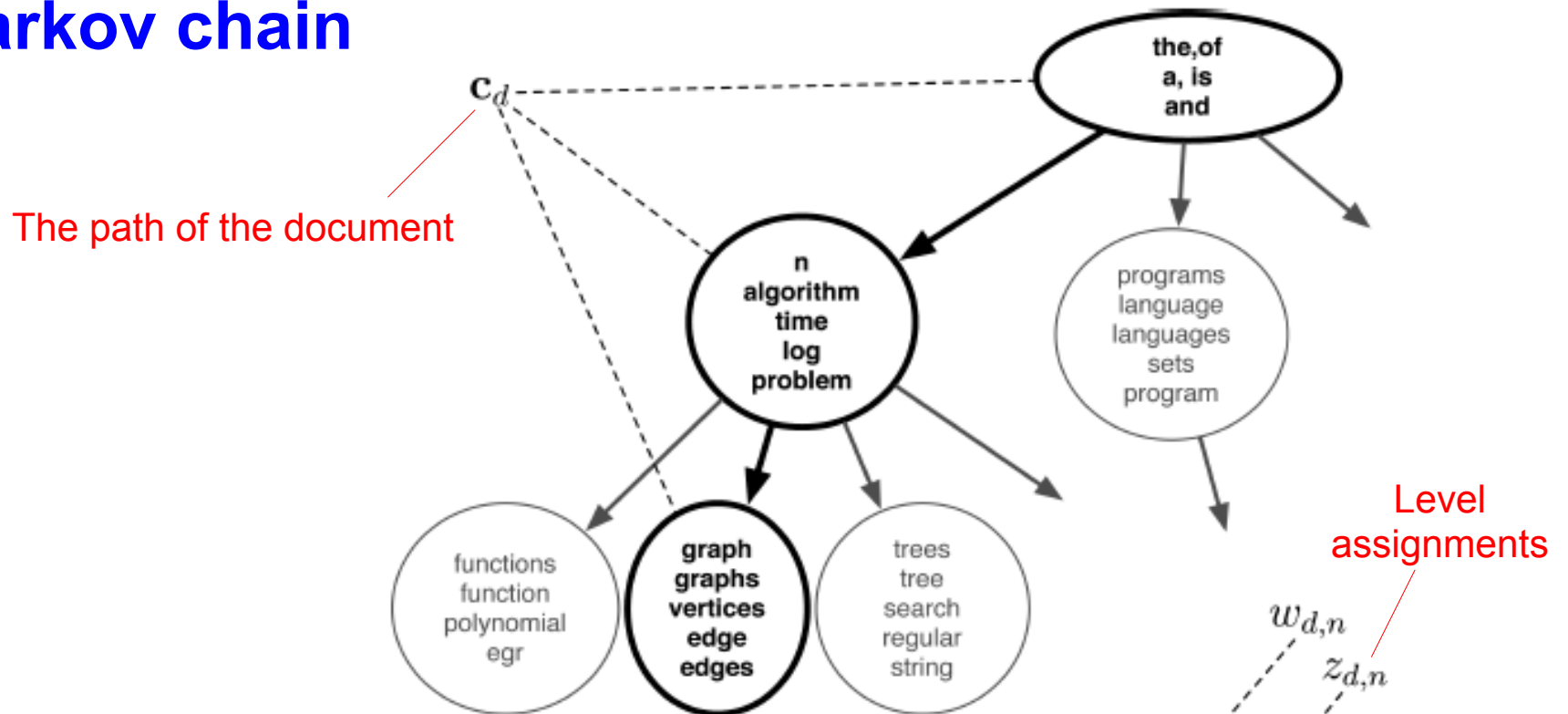
Inference (CGS)

Integrate out distributions over levels and topics

Sample the paths and level assignments



The Markov chain



All₀ previously₁ known₀ efficient₂ maximum-flow algorithms₁ work₁ by₀ finding₁ augmenting paths₁, either₁ one₀ path₀ at₀ a₀ time₁ (as₀ in₀ the₀ original₂ Ford and Fulkerson algorithm₁) or₀ all₀ shortest-length augmenting paths₁ at₀ once₀ (using₀ the₀ layered network₂ approach₁ of₀ Dinic). An₀ alternative₁ method₀ based₀ on₀ the₀ preflow concept₀ of₀ Karzanov is₀ introduced₀. A₀ preflow is₀ like₀ a₀ flow₂ except₁ that₀ the₀ total₀ amount₀ flowing into₁ a₀ vertex₂ is₀ allowed₀ to₀ exceed the₀ total₀ amount₀ flowing out₂. The₀ method₀ maintains a₀ preflow in₀ the₀ original₂ network₂ and₀ pushes local₀ flow₂ excess toward₁ the₀ sink along₀ what₀ are₀ estimated to₀ be₀ shortest₂ paths₁. The₀ algorithm₁ and₀ its₀ analysis₀ are₀ simple₁ and₀ intuitive₁, yet₀ the₀ algorithm₁ runs₀ as₀ fast₁ as₀ any₀ other₀ known₀ method₀ on₀ dense graphs₂ achieving an₀ $O(n)$ time₁ bound₁ on₀ an₀ n -vertex₂ graph₂ by₀ incorporating the₀ dynamic₁ tree₁ data₀ structure₁ of₀ Sleator and₀ Tarjan...

Sampling level assignments

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{c}, \mathbf{w}, m, \pi, \eta) \propto p(z_{d,n} | \mathbf{z}_{d,-n}, m, \pi) p(w_{d,n} | \mathbf{z}, \mathbf{c}, \mathbf{w}_{-(d,n)}, \eta)$$

Probability of level assignment
given other level assignments

Probability of word given
level assignments

Want to sample:

$$p(z_{d,n} = k | \mathbf{z}_{-(d,n)}, \mathbf{c}, \mathbf{w}, m, \pi, \eta), \quad k = 1, 2, \dots$$

Infinite

Solution: sample in stages

- Sample a level from all currently represented levels, plus one level deeper
- If deeper level is chosen, sample from a Bernoulli to go even deeper, iteratively

Sampling level assignments

For already represented levels:

$$\begin{aligned}
 p(z_{d,n} = k | \mathbf{z}_{d,-n}, m, \pi) &= \mathbb{E} \left[V_k \prod_{j=1}^{k-1} (1 - V_j) \right] \quad \leftarrow \text{Expected length of the } k\text{-th stick} \\
 &= \mathbb{E}[V_k] \prod_{j=1}^{k-1} \mathbb{E}[1 - V_j] \\
 &= \frac{m\pi + \#[\mathbf{z}_{d,-n} = k]}{\pi + \#[\mathbf{z}_{d,-n} \geq k]} \prod_{j=1}^{k-1} \frac{(1 - m)\pi + \#[\mathbf{z}_{d,-n} > j]}{\pi + \#[\mathbf{z}_{d,-n} \geq j]}
 \end{aligned}$$

$$p(w_{d,n} | \mathbf{z}, \mathbf{c}, \mathbf{w}_{-(d,n)}, \eta) \propto \#[\mathbf{z}_{-(d,n)} = z_{d,n}, \mathbf{c}_{z_{d,n}} = \mathbf{c}_{d,z_{d,n}}, \mathbf{w}_{-(d,n)} = w_{d,n}] + \eta$$

Smoothed number of similar words assigned to this level in this path

For the deeper levels:

$$p(z_{d,n} > \max(\mathbf{z}_{d,-n}) | \mathbf{z}_{d,-n}, \mathbf{w}, m, \pi, \eta) = 1 - \sum_{j=1}^{\max(\mathbf{z}_{d,-n})} p(z_{d,n} = j | \mathbf{z}_{d,-n}, \mathbf{w}, m, \pi, \eta)$$

Leftover probability after we account for the already represented levels

Sampling level assignments

If we choose to go deeper...

Start with:

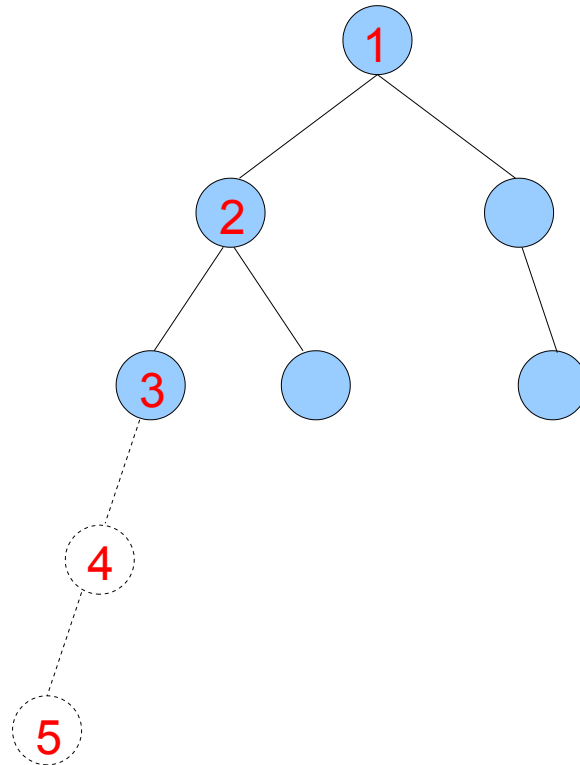
$$\ell = \max(\mathbf{z}_{d,-n}) + 1$$

And sample from a Bernoulli with parameters:

$$p(z_{d,n} = \ell \mid z_{d,-n}, z_{d,n} > \ell - 1, \mathbf{w}, m, \pi, \eta) = (1 - m)p(w_{d,n} \mid \mathbf{z}, \mathbf{c}, \mathbf{w}_{-(d,n)}, \eta)$$

$$p(z_{d,n} > \ell \mid z_{d,-n}, z_{d,n} > \ell - 1) = 1 - p(z_{d,n} = \ell \mid z_{d,-n}, z_{d,n} > \ell - 1, \mathbf{w}, m, \pi, \eta).$$

Until we get a “hit”.



Sampling paths

We're only concerned with paths of the length of the deepest level assignment variable for the document under consideration.

$$p(\mathbf{c}_d | \mathbf{w}, \mathbf{c}_{-d}, \mathbf{z}, \eta, \gamma) \propto p(\mathbf{c}_d | \mathbf{c}_{-d}, \gamma) p(\mathbf{w}_d | \mathbf{c}, \mathbf{w}_{-d}, \mathbf{z}, \eta)$$

Probability of path given
all other paths

Probability of words
given path

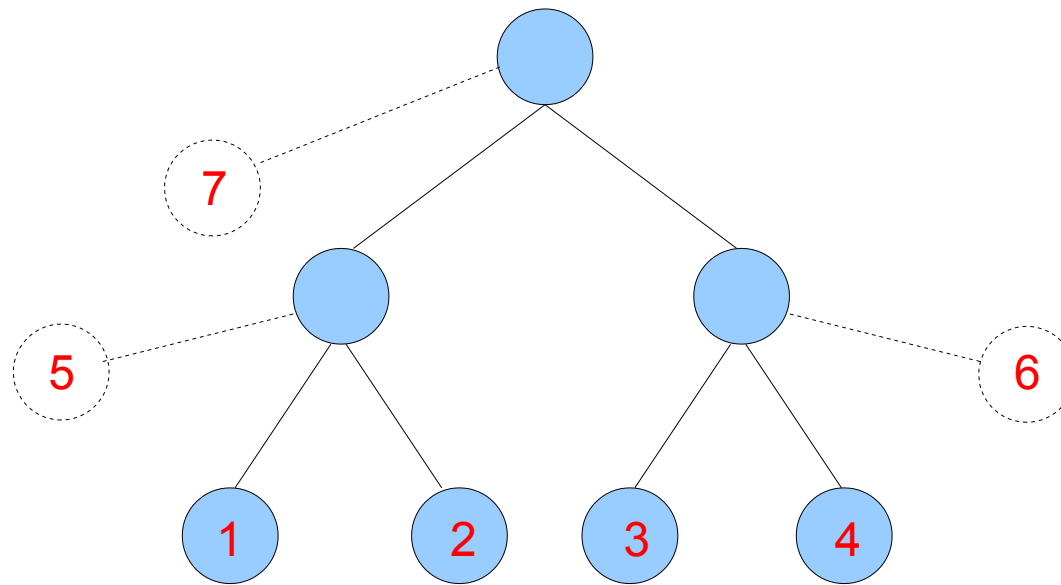
$$p(\mathbf{c}_d | \mathbf{c}_{-d}, \gamma) = \text{nCRP prior}$$

$$p(\mathbf{w}_d | \mathbf{c}, \mathbf{w}_{-d}, \mathbf{z}, \eta) = \prod_{\ell=1}^{\max(\mathbf{z}_d)} \frac{\Gamma(\sum_w \#[\mathbf{z}_d = \ell, \mathbf{c}_{-d, \ell} = c_{d, \ell}, \mathbf{w}_{-d} = w] + V\eta)}{\prod_w \Gamma(\#[\mathbf{z}_{-d} = \ell, \mathbf{c}_{-d, \ell} = c_{d, \ell}, \mathbf{w}_{-d} = w] + \eta)} \times \frac{\prod_w \Gamma(\#[\mathbf{z} = \ell, \mathbf{c}_\ell = c_{d, \ell}, \mathbf{w} = w] + \eta)}{\Gamma(\sum_w \#[\mathbf{z} = \ell, \mathbf{c}_\ell = c_{d, \ell}, \mathbf{w} = w] + V\eta)}$$

Note: path must be drawn as a block; levels are not independent.

Sampling paths

All the possible paths a new document can take through the existing tree:



Sampling the hyperparameters

Take the nice Bayesian approach of putting priors on the hyperparameters:

$$m \sim \text{Beta}(\alpha_1, \alpha_2)$$

$$\pi \sim \text{Exponential}(\alpha_3)$$

$$\gamma \sim \text{Gamma}(\alpha_4, \alpha_5)$$

$$\eta \sim \text{Exponential}(\alpha_6)$$

Put Metropolis-Hastings steps between iterations of the Gibbs sampler.

What purpose does this serve?

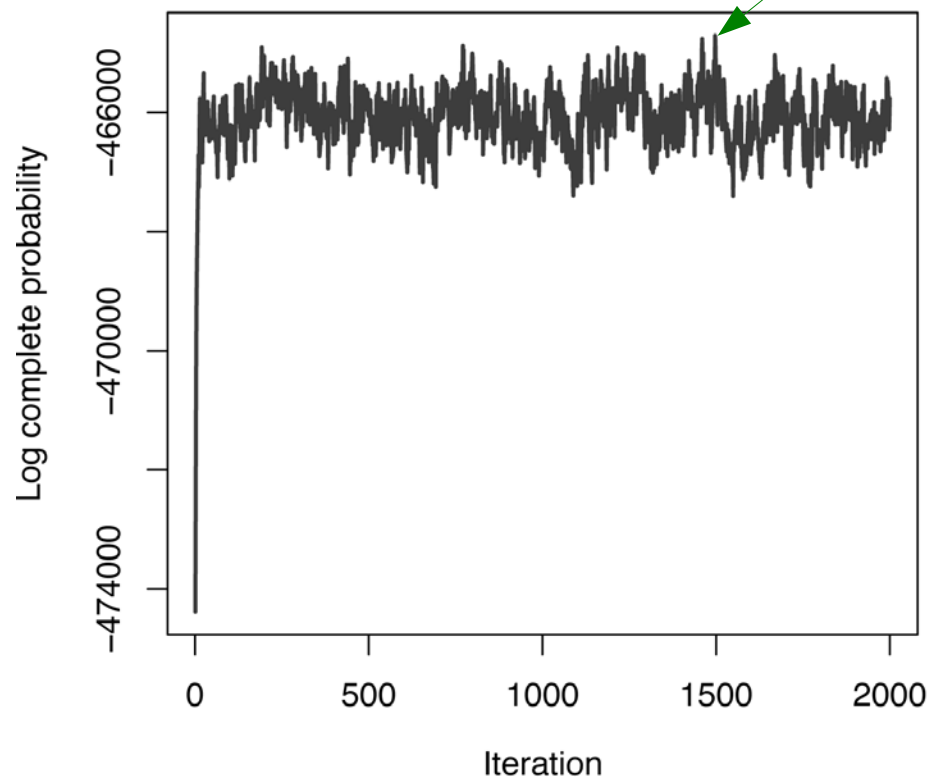
Resulting inference is less influenced by hyper-hyperparameters than by hyperparameters.

Convergence

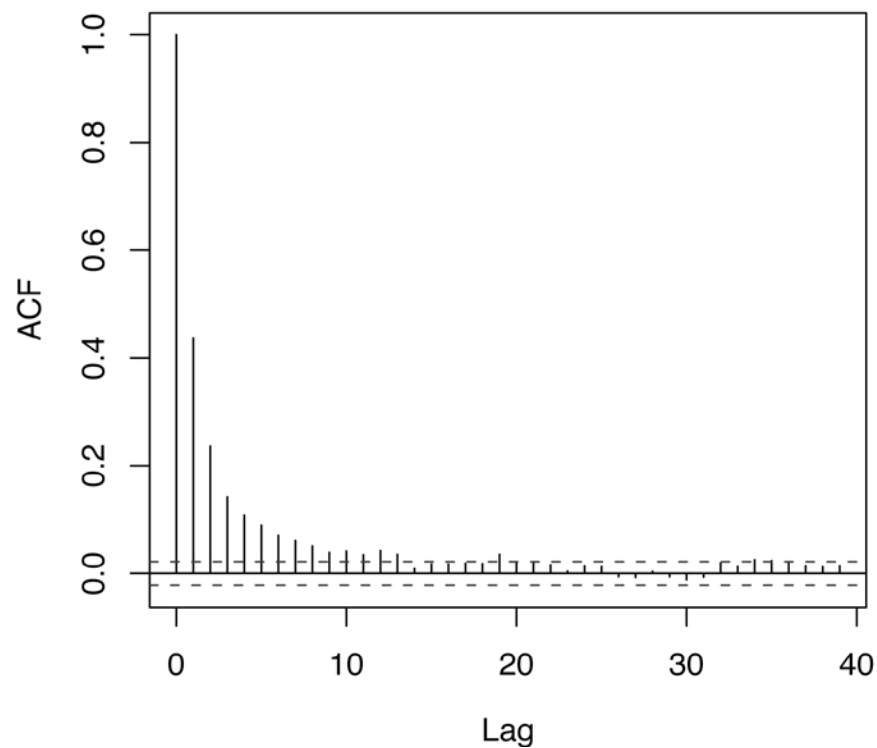
$$\mathcal{L}^{(t)} = \log p \left(\mathbf{c}_{1:D}^{(t)}, \mathbf{z}_{1:D}^{(t)}, \mathbf{w}_{1:D} | \gamma, \eta, m, \pi \right)$$

Approximation of the posterior mode for these iterations

Log complete probability for the JACM corpus



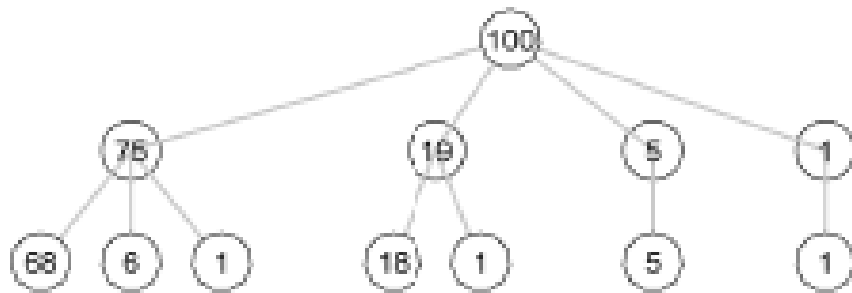
Autocorrelation function for the JACM corpus



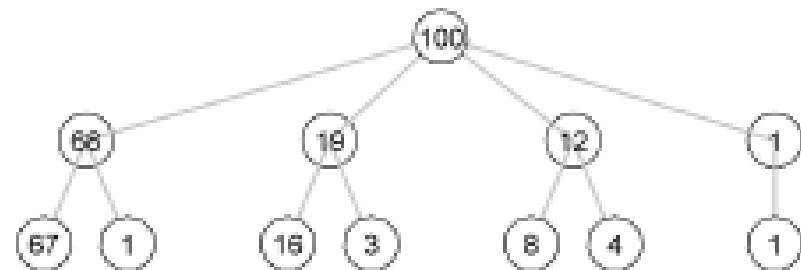
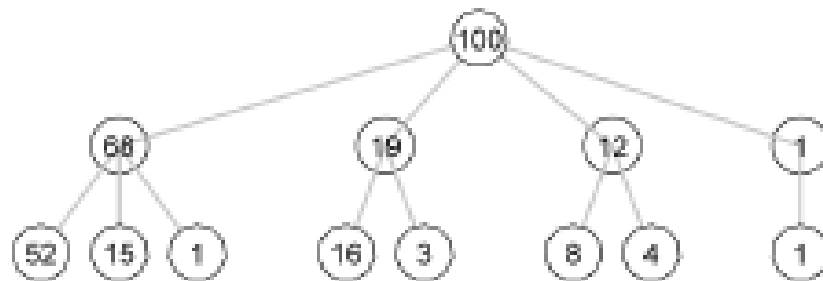
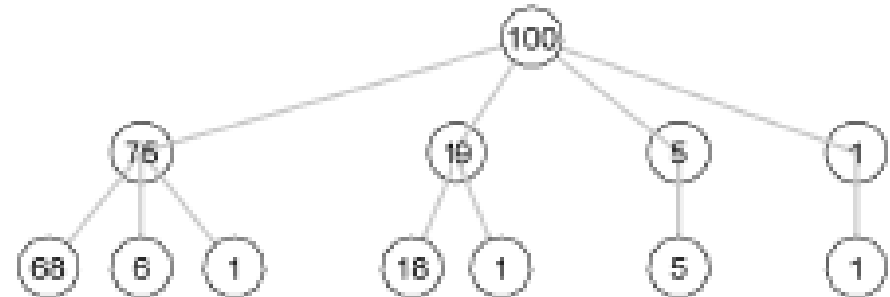
Experiments (simulated data)

- 100 documents drawn from an hLDA model.
- $\eta = 0.005$; $\gamma = 1$; $V = 100$

True dataset hierarchy



Posterior mode



Experiments (scientific abstracts)

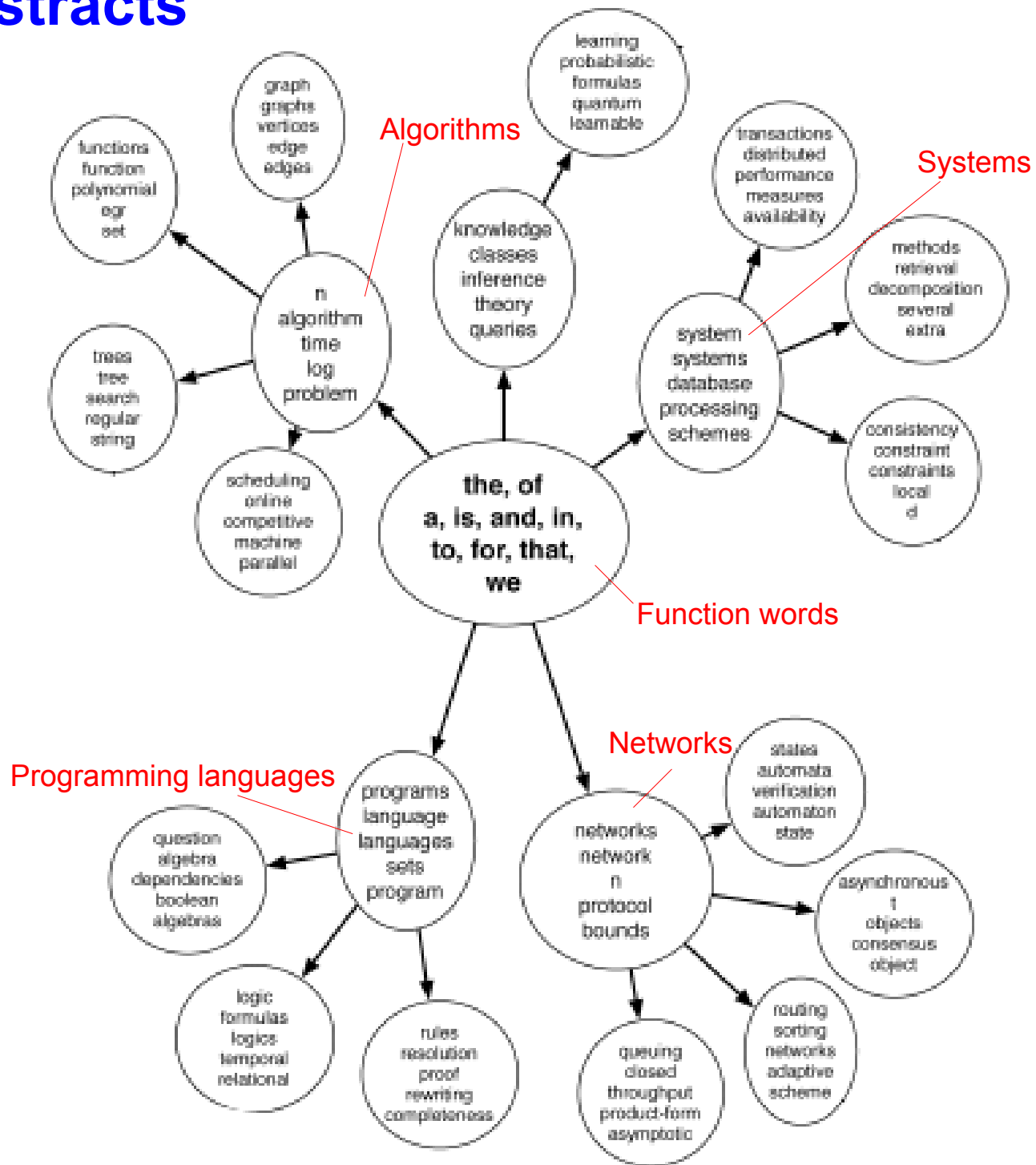
- Fix parameters:

$$\eta = \{2.0, 1.0, 0.5\}; \quad \gamma = 1.0; \quad \pi = 100; \quad m = 0.5$$

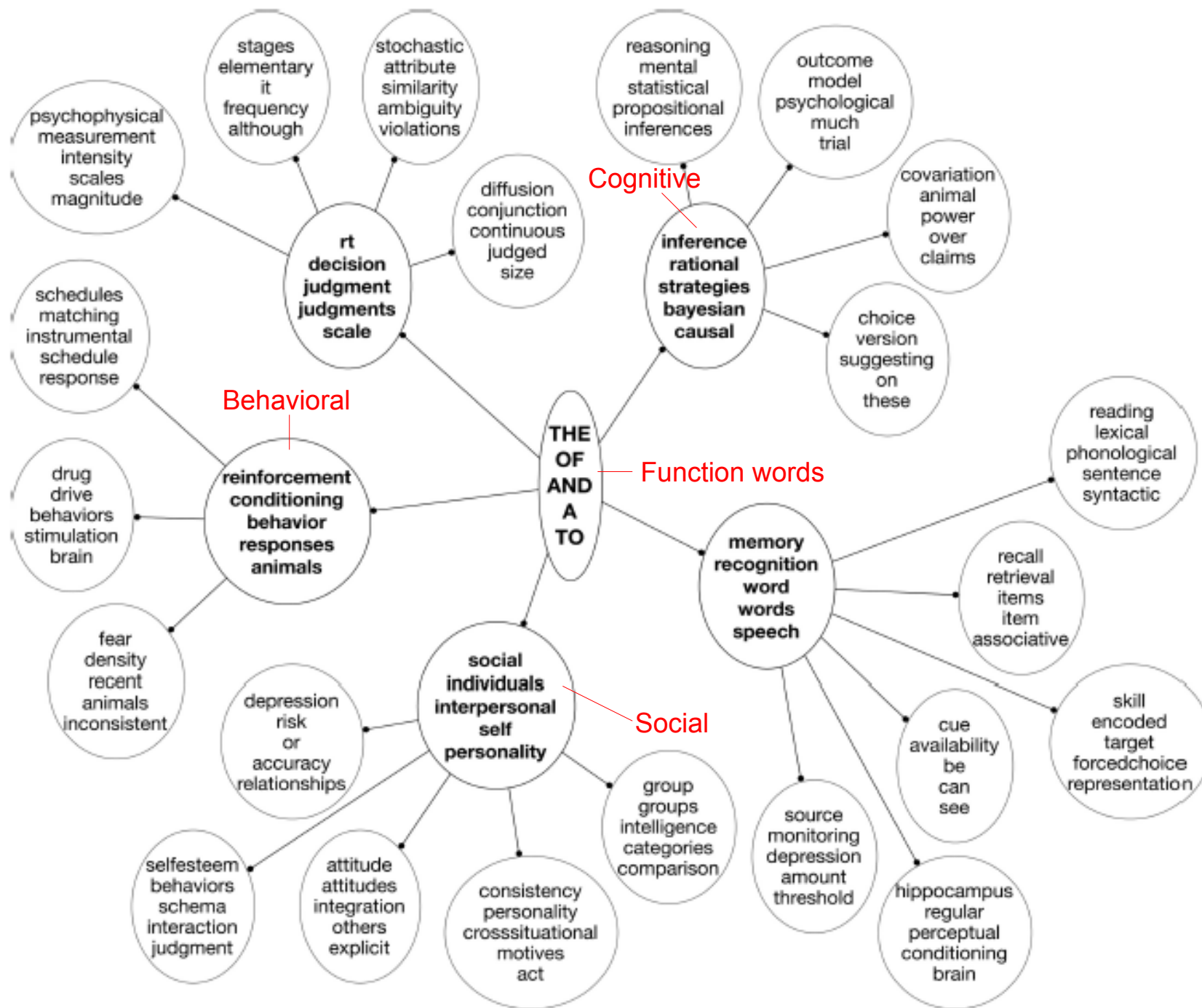
- Recall that the topics have been integrated out, so they'll have to be estimated in the posterior. Posterior inference only yields a tree structure \mathbf{c} , and assignments to levels \mathbf{z} .
- The probability of a particular word w at a particular level l in a particular path \mathbf{p} is:

$$p(w|\mathbf{z}, \mathbf{c}, \mathbf{w}, \eta) = \frac{\#[\mathbf{z} = \ell, \mathbf{c} = \mathbf{p}, \mathbf{w} = w] + \eta}{\#[\mathbf{z} = \ell, \mathbf{c} = \mathbf{p}] + V\eta}$$

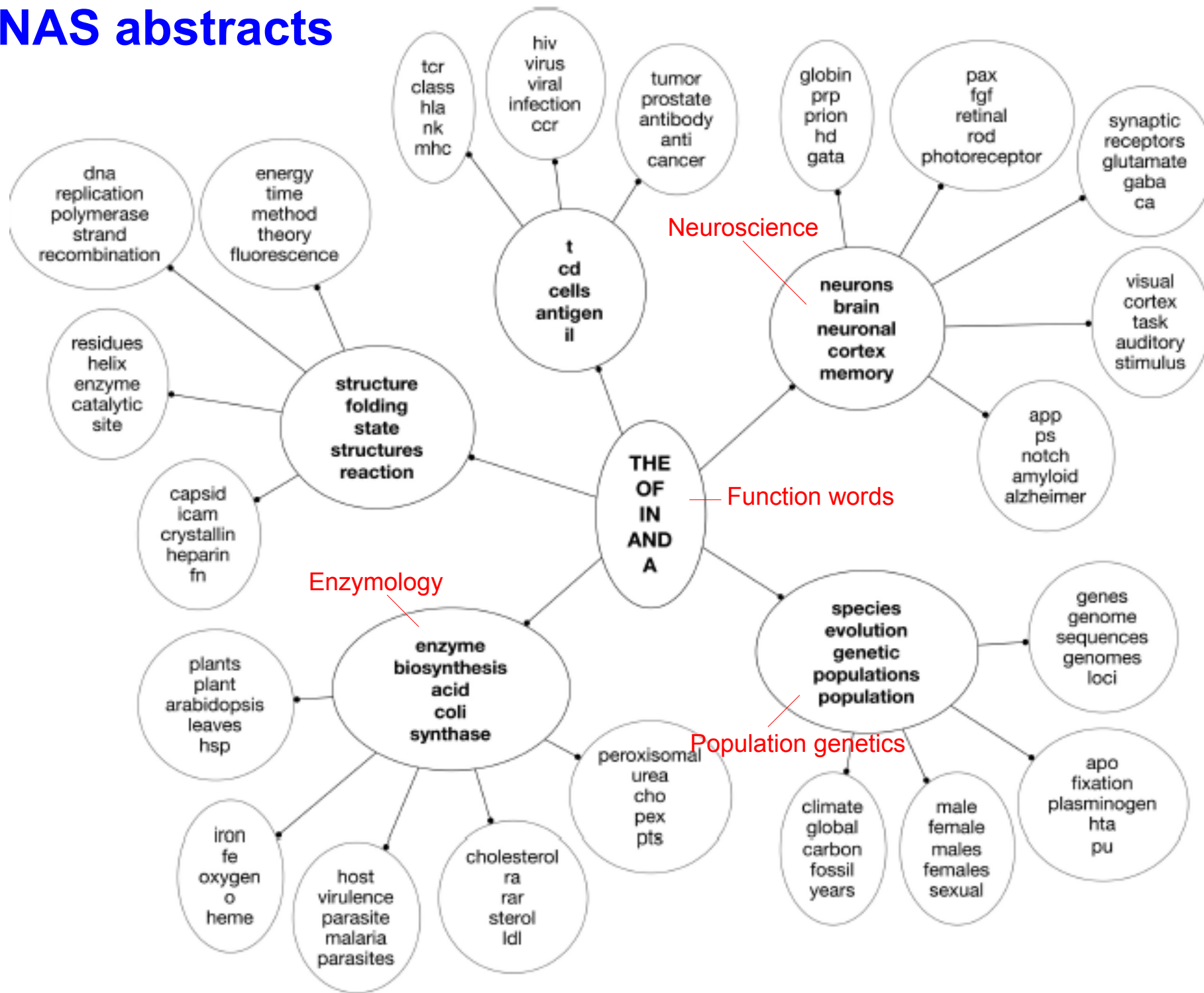
JACM abstracts



Psychological Review abstracts

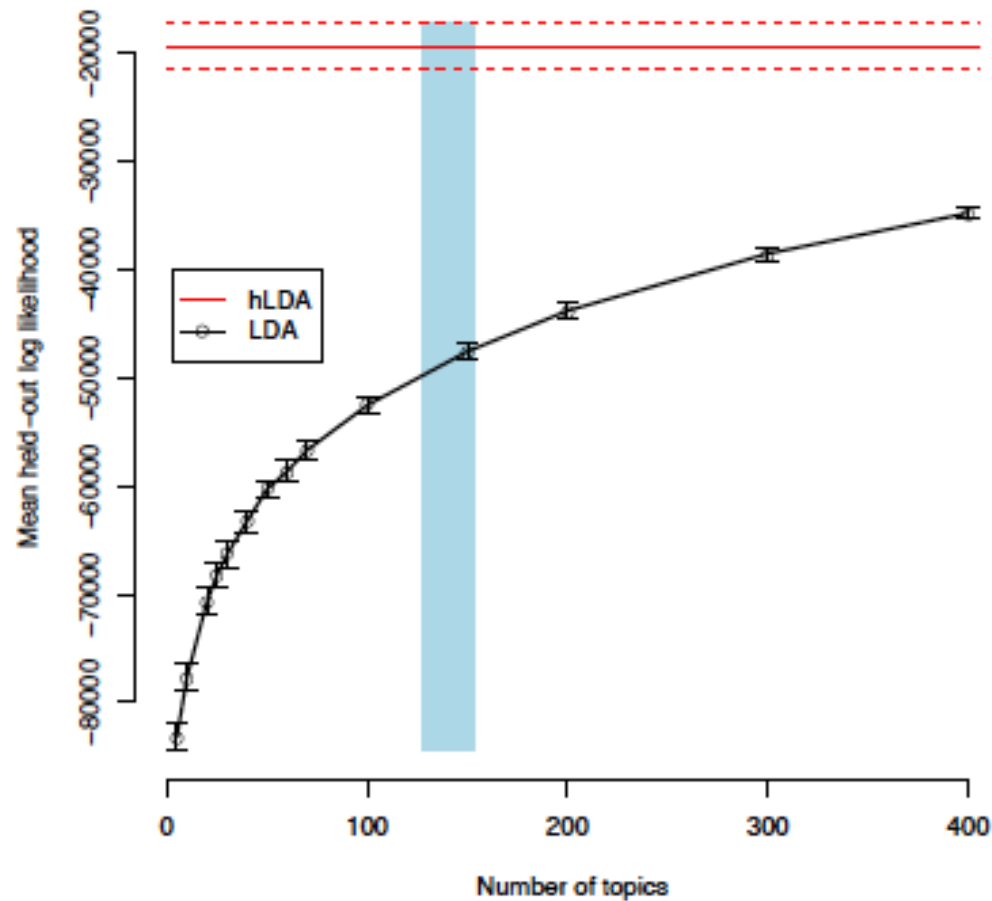


PNAS abstracts

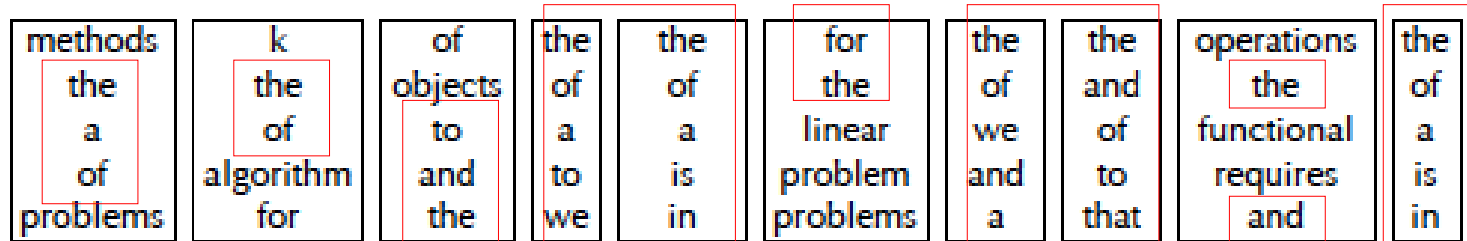


LDA v. hLDA

- Number of topics in LDA fixed beforehand
- In LDA, each document can place an arbitrary distribution over topics, not only those that lie on a path in the hierarchy.



LDA v. hLDA



Function words everywhere

Conclusions

- **The Good:**

- The nested CRP allows a flexible family of prior distributions over arbitrary tree structures; definitely could be useful for more than just topic models.
- Nice qualitative results for topic hierarchies.
- Same inference of number of topics as a model like HDP

- **The Bad/Ugly:**

- The restriction that documents can only follow a single path in the tree is a possibly limiting one. (Michael Jordan talked about this extension when he spoke in class).
- Quantitative evaluation is not extensive enough.
- I'd like to see comparisons of hLDA with HDP, as opposed to LDA. It seems like that would get closer to the heart of whether hierarchies are helpful or not.