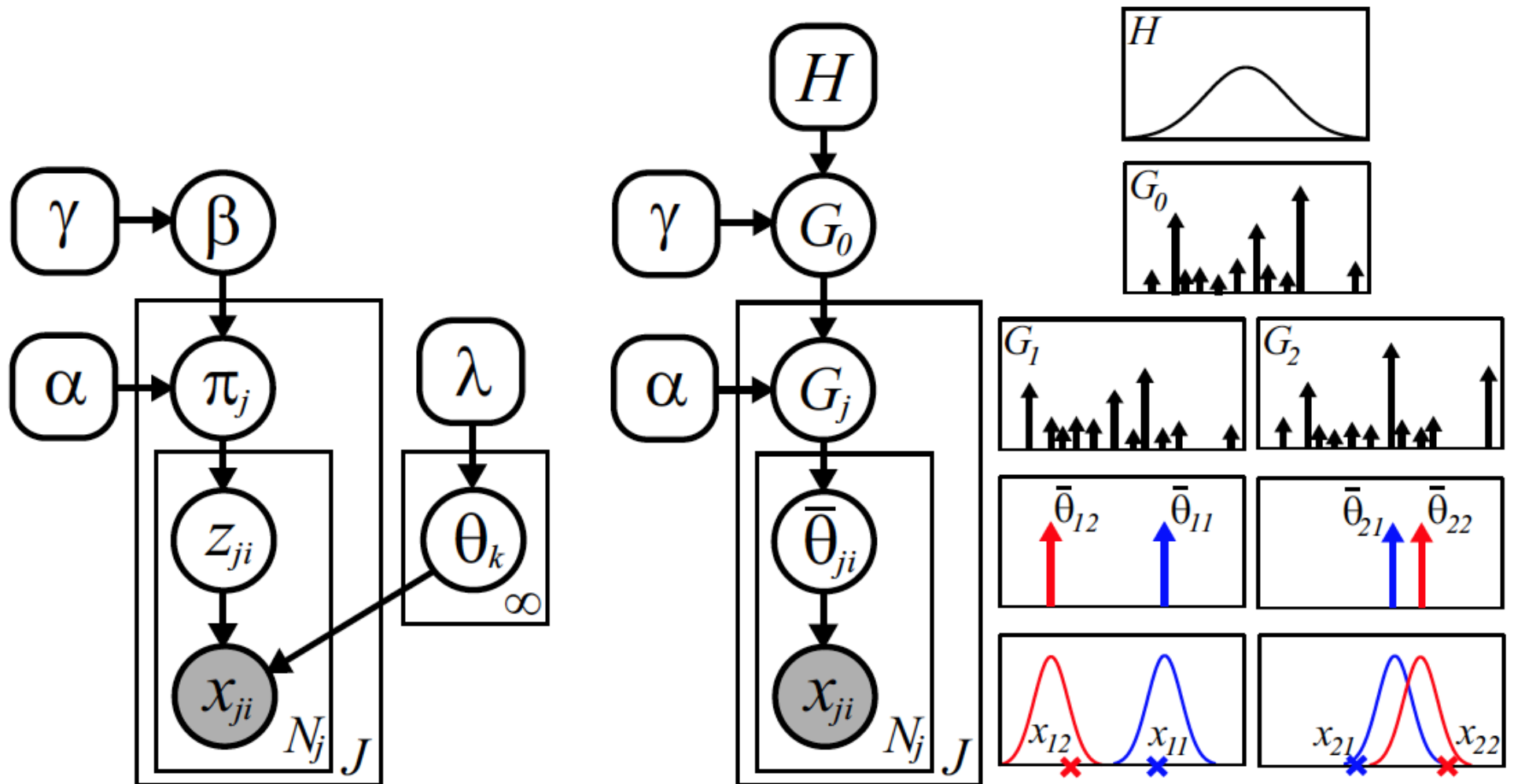


Applied Bayesian Nonparametrics

Special Topics in Machine Learning
Brown University CSCI 2950-P, Fall 2011

October 27: Pitman-Yor Processes,
Infinite Markov Models

Hierarchical Dirichlet Process



Hierarchical Dirichlet Process

$$G_0(\theta) = \sum_{k=1}^{\infty} \beta_k \delta(\theta, \theta_k)$$

$$\beta \sim \text{GEM}(\gamma)$$

$$\theta_k \sim H(\lambda) \quad k = 1, 2, \dots$$

$$G_j(\theta) = \sum_{t=1}^{\infty} \tilde{\pi}_{jt} \delta(\theta, \tilde{\theta}_{jt})$$

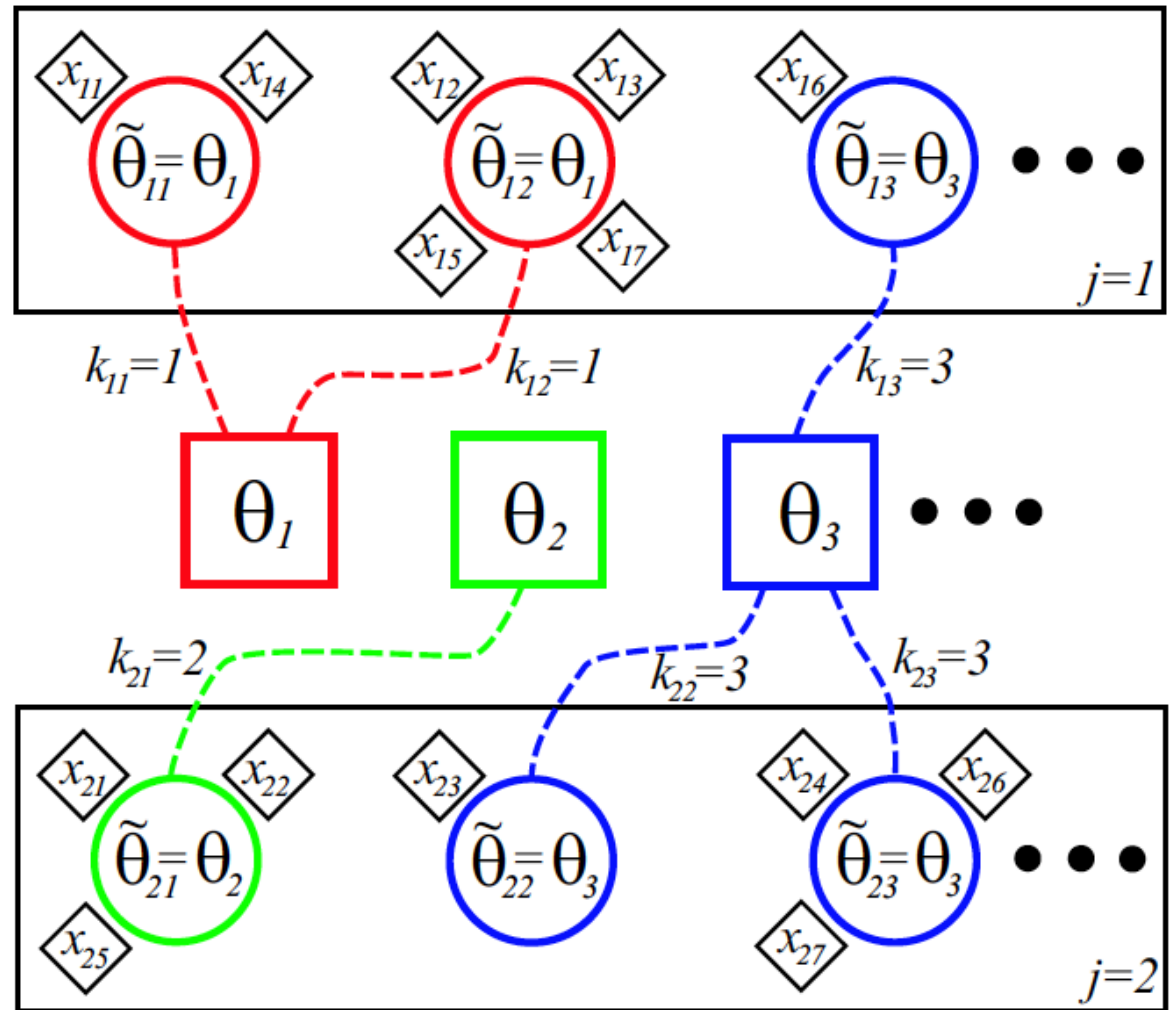
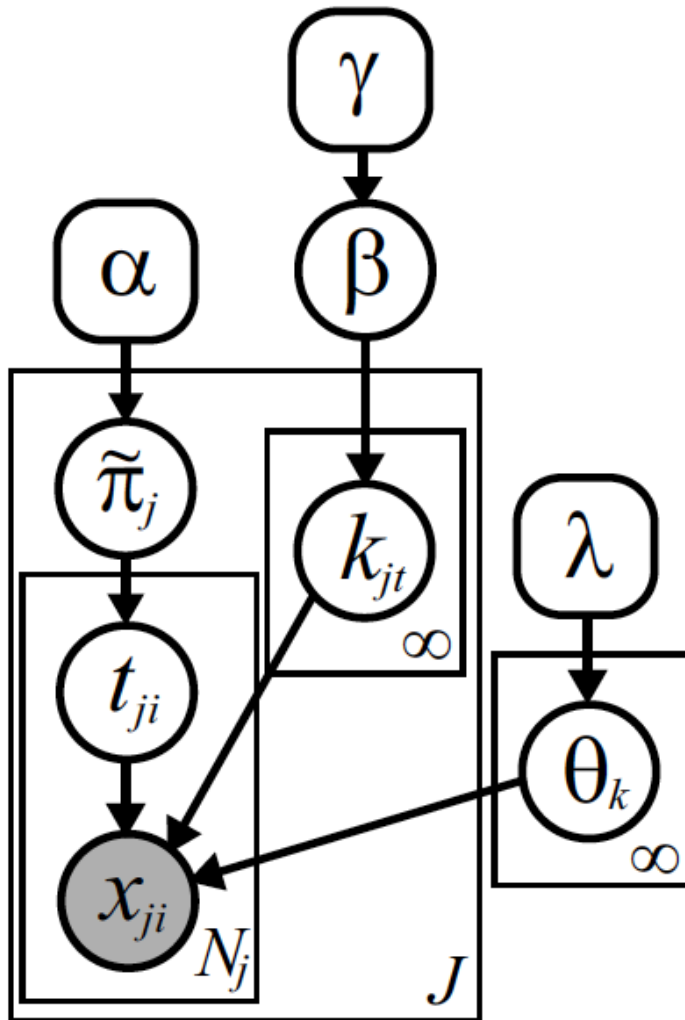
$$\tilde{\pi}_j \sim \text{GEM}(\alpha)$$

$$\tilde{\theta}_{jt} \sim G_0 \quad t = 1, 2, \dots$$

$$G_j(\theta) = \sum_{k=1}^{\infty} \pi_{jk} \delta(\theta, \theta_k)$$

$$\pi_{jk} = \sum_{t|k_{jt}=k} \tilde{\pi}_{jt}$$

Chinese Restaurant Franchise



$$p(t_{ji} | t_{j1}, \dots, t_{ji-1}, \alpha) \propto \sum_t N_{jt} \delta(t_{ji}, t) + \alpha \delta(t_{ji}, \bar{t})$$

$$p(k_{jt} | \mathbf{k}_1, \dots, \mathbf{k}_{j-1}, k_{j1}, \dots, k_{jt-1}, \gamma) \propto \sum_k M_k \delta(k_{jt}, k) + \gamma \delta(k_{jt}, \bar{k})$$

Views of Hierarchical Dirichlet

Construction repeated at each level of hierarchy:

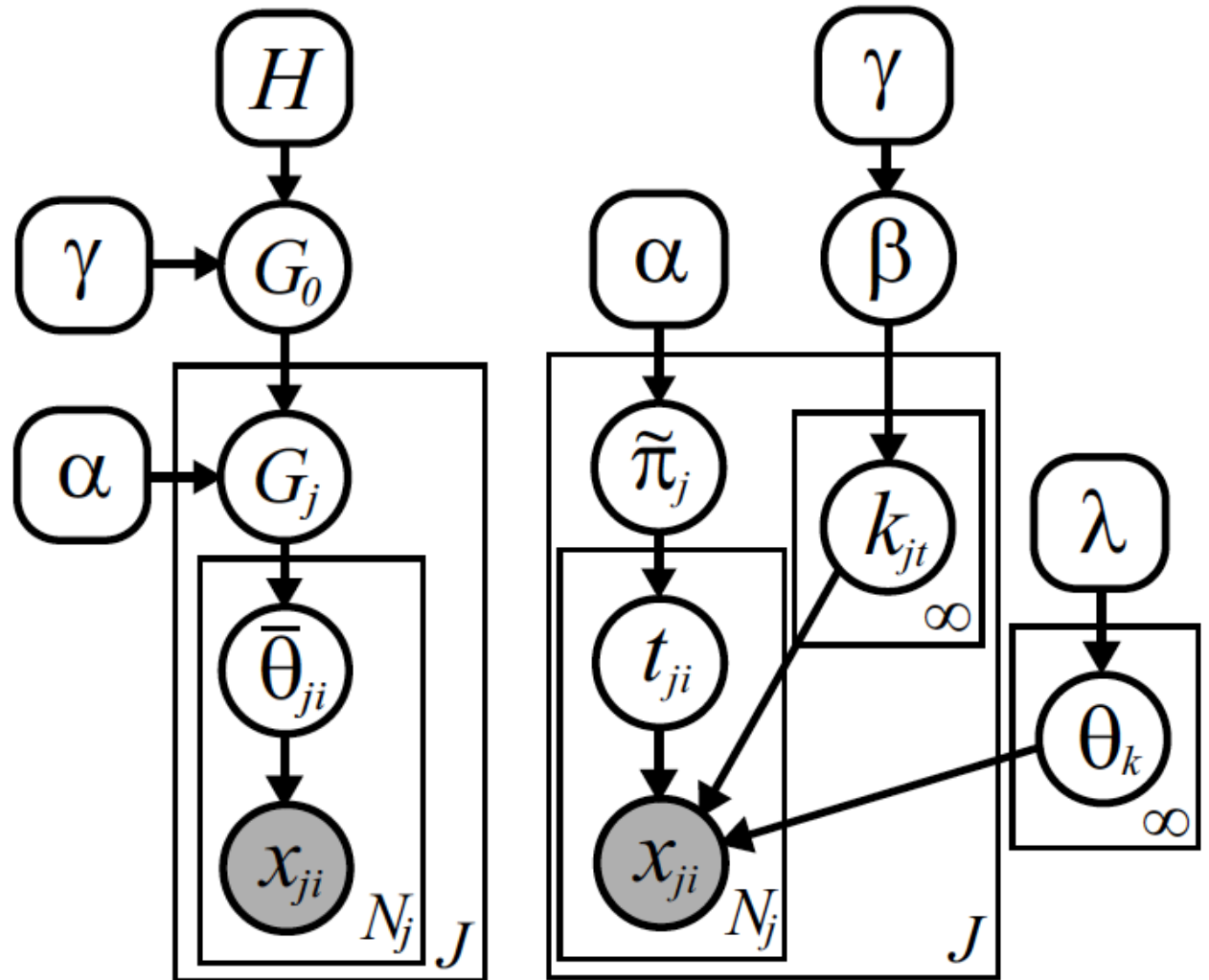
$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}$$

$$\beta_k = v_k \prod_{\ell=1}^{k-1} (1 - v_\ell)$$

$$v_k \sim \text{Beta}(1, \gamma)$$

For some ordering of all tables in all groups:

$$p(k_{N+1} \mid k_N, \dots, k_1, \gamma) = \frac{1}{N + \gamma} \left[\sum_{k=1}^K M_k \delta(k_{N+1}, k) + \gamma \delta(k_{N+1}, K + 1) \right]$$



Views of Hierarchical Pitman-Yor

Construction repeated at each level of hierarchy:

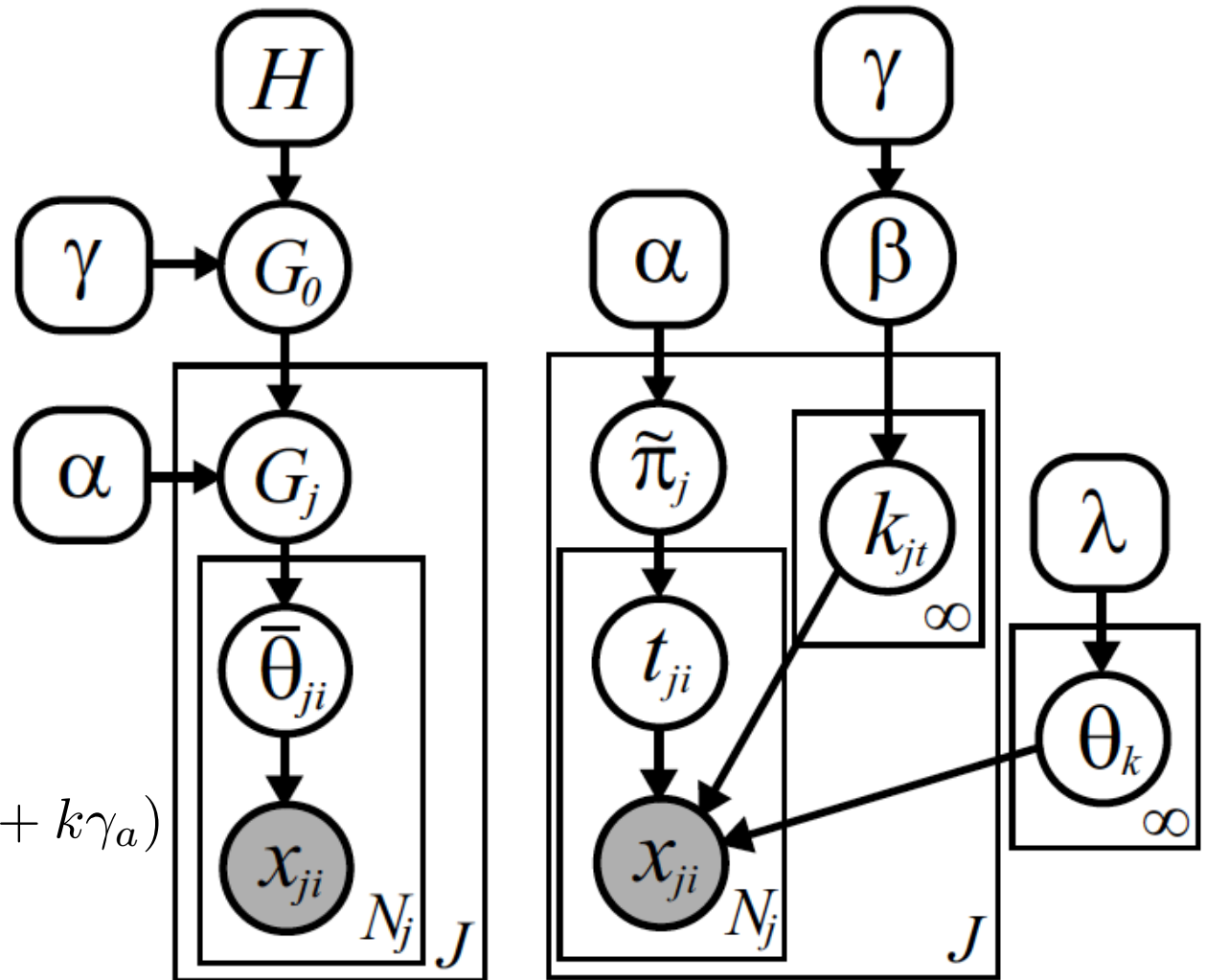
$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}$$

$$\beta_k = v_k \prod_{\ell=1}^{k-1} (1 - v_\ell)$$

$$v_k \sim \text{Beta}(1 - \gamma_a, \gamma_b + k\gamma_a)$$

For some ordering of all tables in all groups:

$$p(k_{N+1} | k_N, \dots, k_1, \gamma) = \frac{1}{N + \gamma_b} \left[\sum_{k=1}^K (M_k - \gamma_a) \delta(k_{N+1}, k) + (\gamma_b + K\gamma_a) \delta(k_{N+1}, K + 1) \right]$$



Why Pitman-Yor?

Generalizing the Dirichlet Process

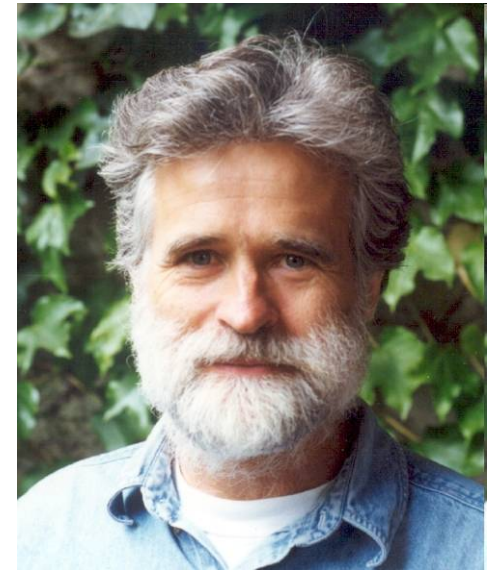
- Distribution on partitions leads to a generalized *Chinese restaurant process*
- Special cases arise as excursion lengths for Markov chains, Brownian motions, ...

Power Law Distributions

	DP	PY
Number of unique clusters in N observations	$\mathcal{O}(b \log N)$	Heaps' Law: $\mathcal{O}(bN^a)$
Size of sorted cluster weight k	$\mathcal{O}\left(\alpha_b \left(\frac{1+b}{b}\right)^{-k}\right)$	Zipf's Law: $\mathcal{O}\left(\alpha_{ab} k^{-\frac{1}{a}}\right)$

**Natural Language
Statistics**

Goldwater, Griffiths, & Johnson, 2005
Teh, 2006

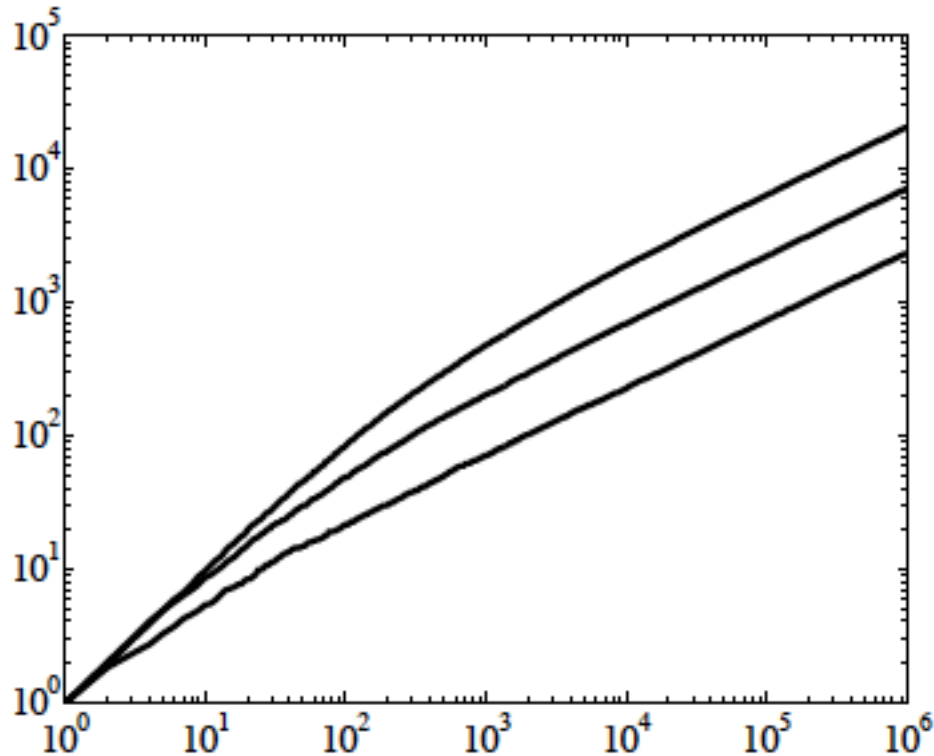


Jim Pitman

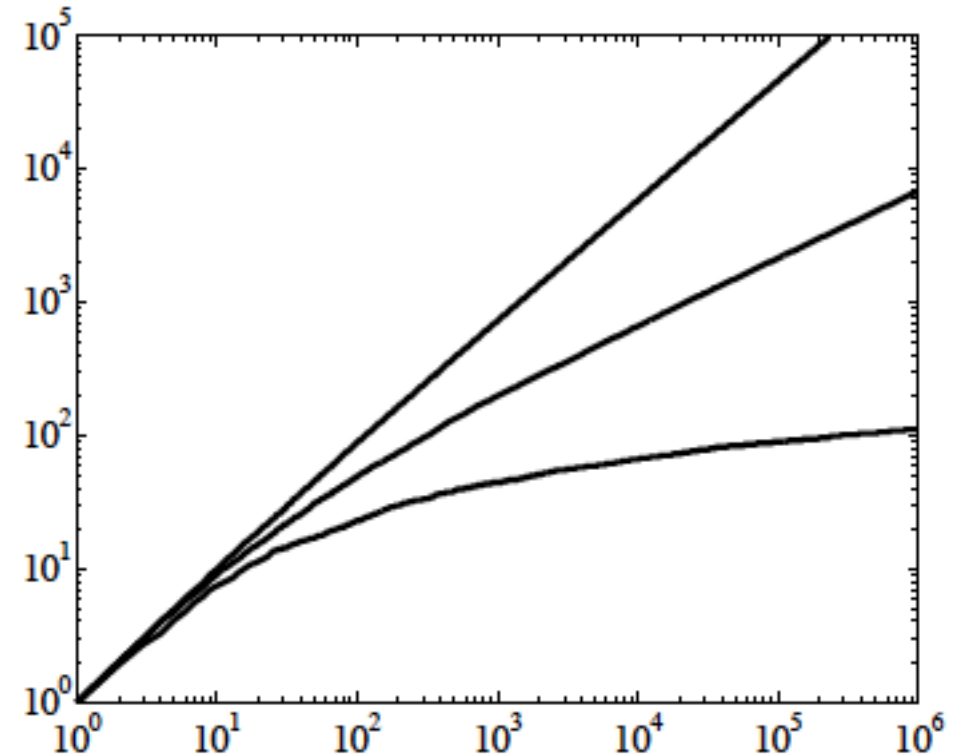


Marc Yor

Number of Unique Words

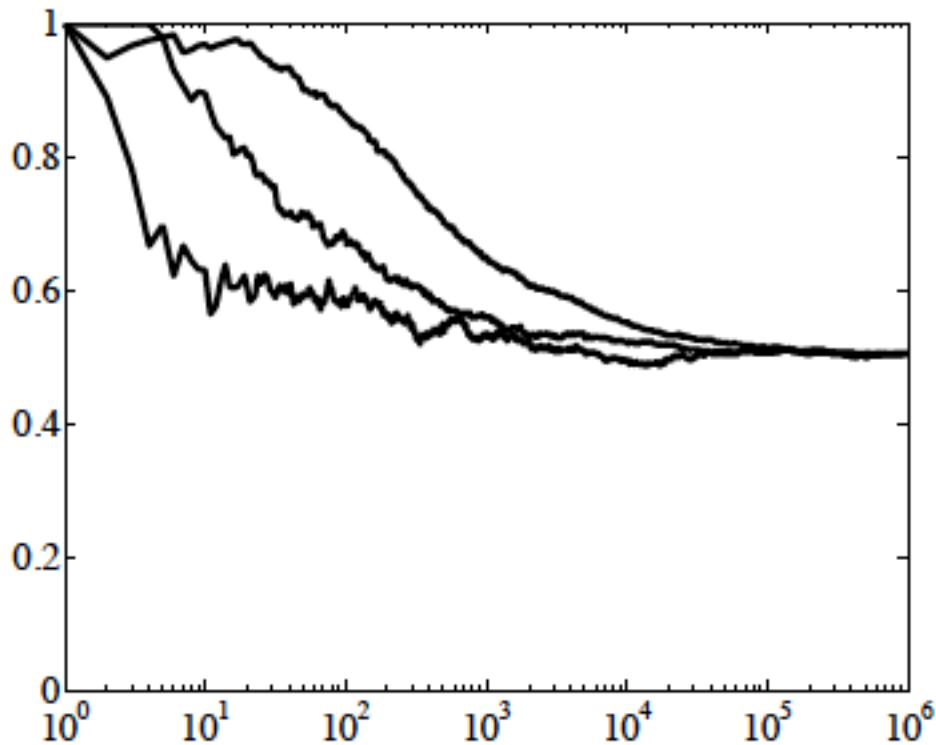


Fixed discount parameter

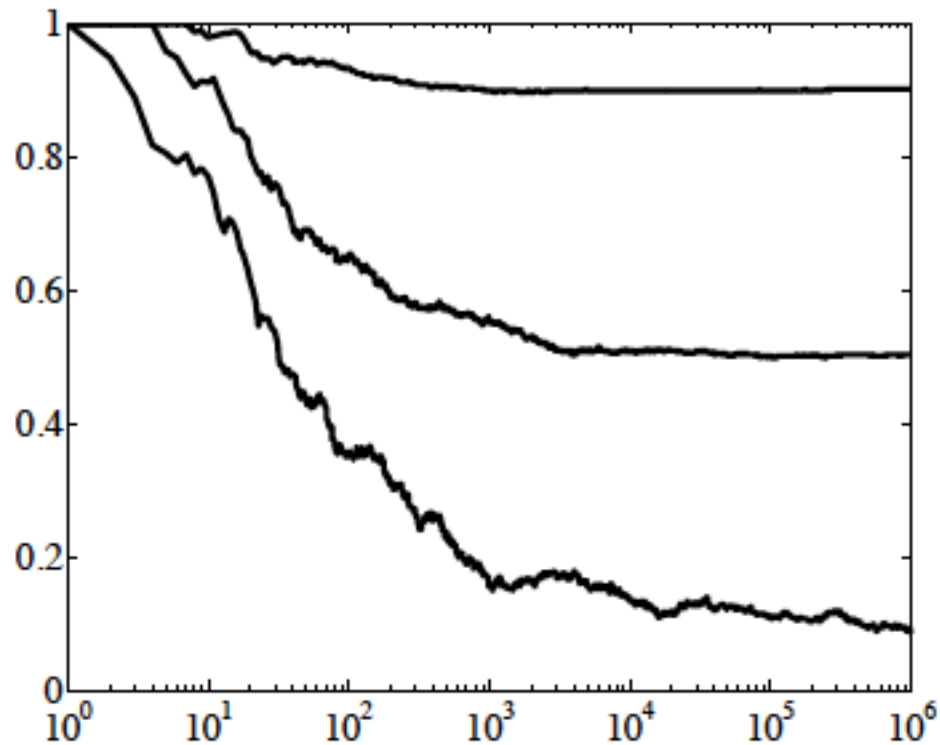


Fixed concentration parameter

Fraction of Words Appearing Once

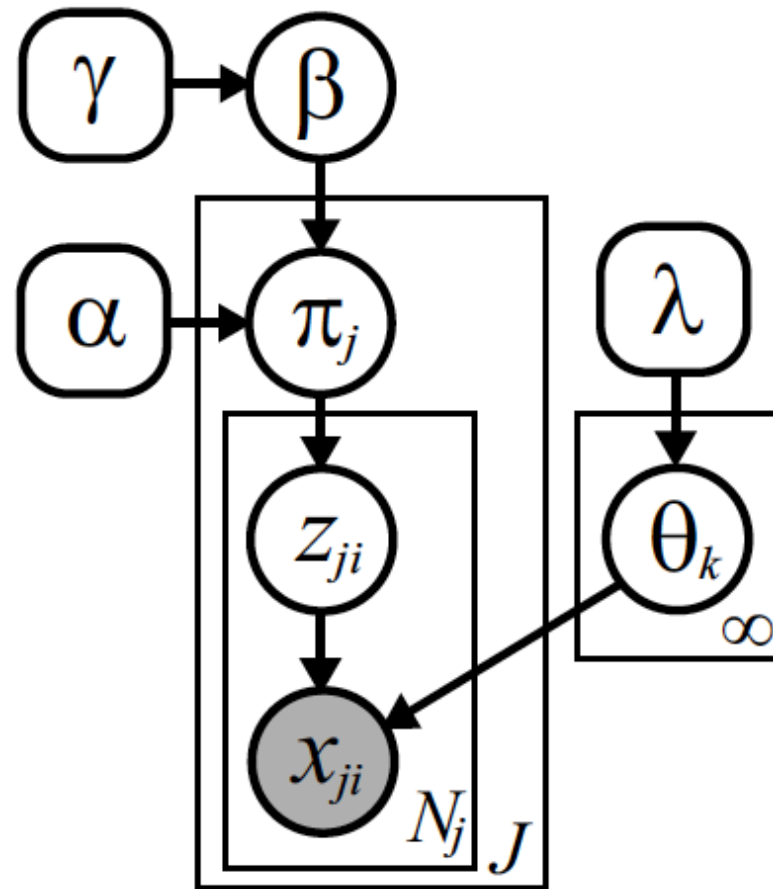


Fixed discount parameter

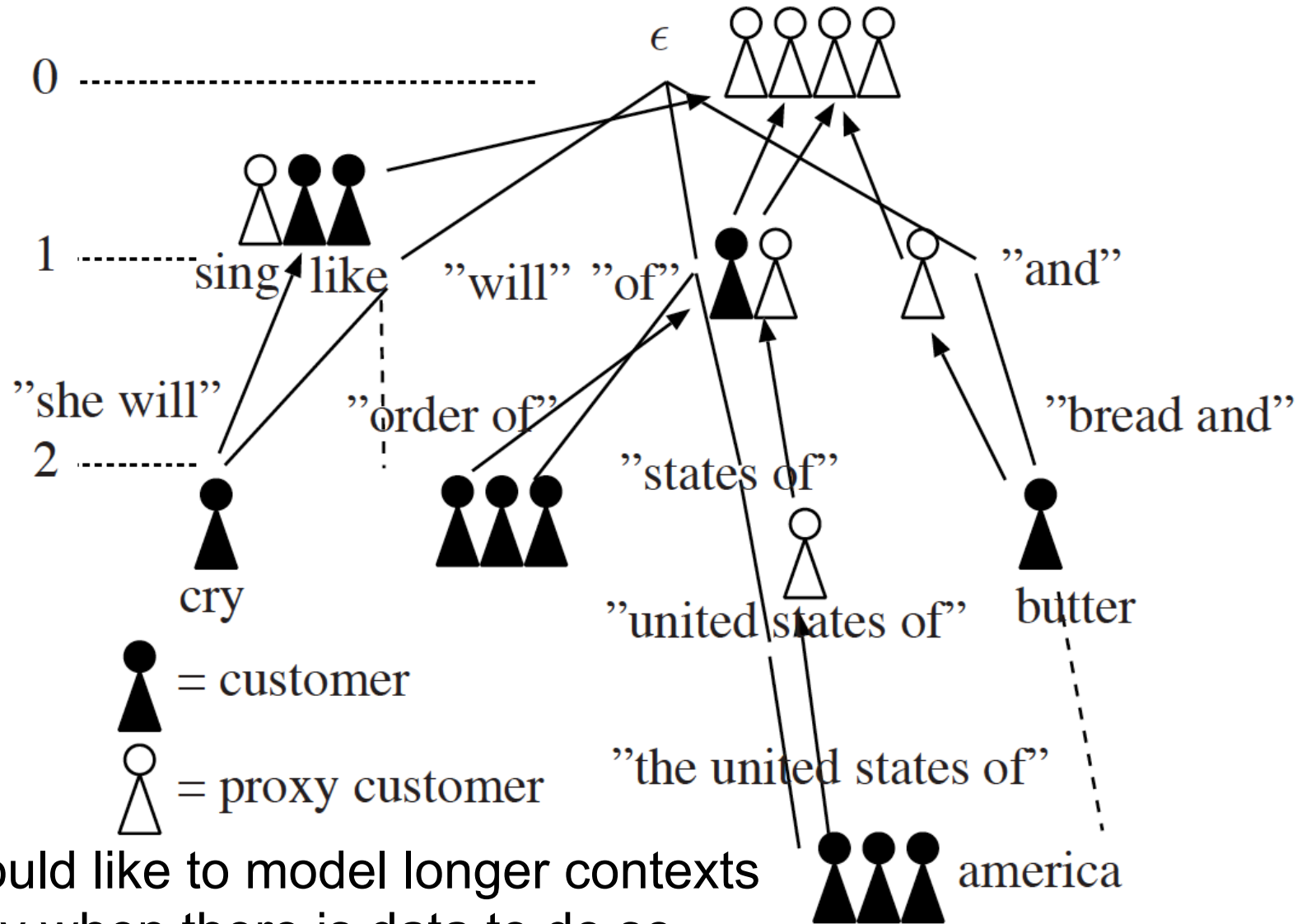


Fixed concentration parameter

A Third View: Tractable for PY?

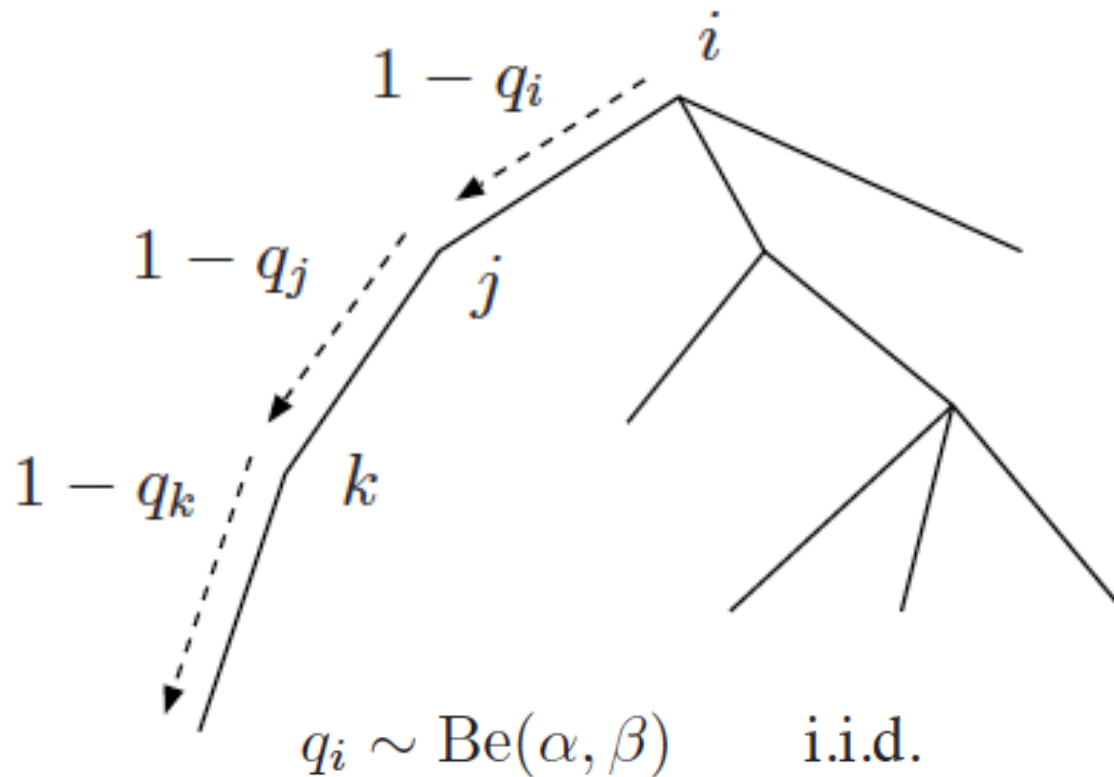


Variable-Order Suffix Tree



- Would like to model longer contexts only when there is data to do so
- How can we make this dynamic?

Making Order Random Under Prior



$$p(n = l|h) = q_l \prod_{i=0}^{l-1} (1 - q_i)$$

$$p(n_t = l | \mathbf{s}_{-t}, \mathbf{z}_{-t}, \mathbf{n}_{-t}) = \frac{a_l + \alpha}{a_l + b_l + \alpha + \beta} \prod_{i=0}^{l-1} \frac{b_i + \beta}{a_i + b_i + \alpha + \beta}$$

Example: Character Model

‘how queershaped little children drawling-desks, which would get through that dormouse!’
said alice; ‘let us all for anything the secondly, but it to have and another question, but i
shalled out, ‘you are old,’ said the you’re trying to far out to sea.

(a) Random walk generation from a character model.

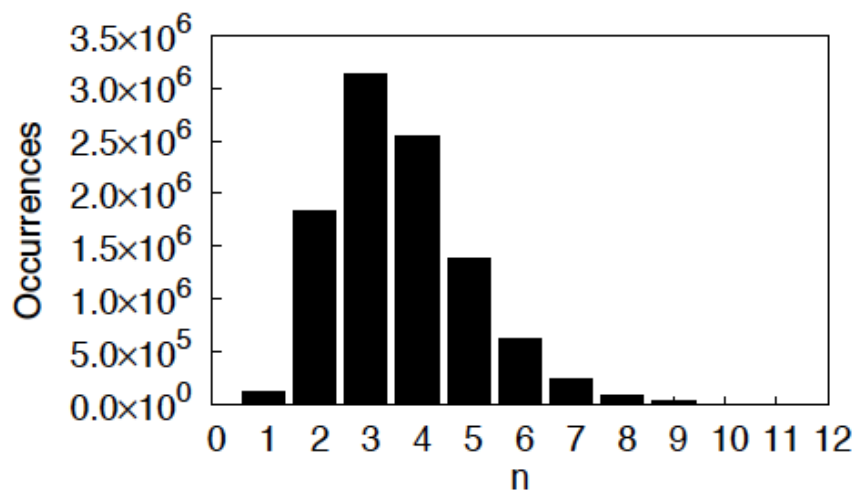
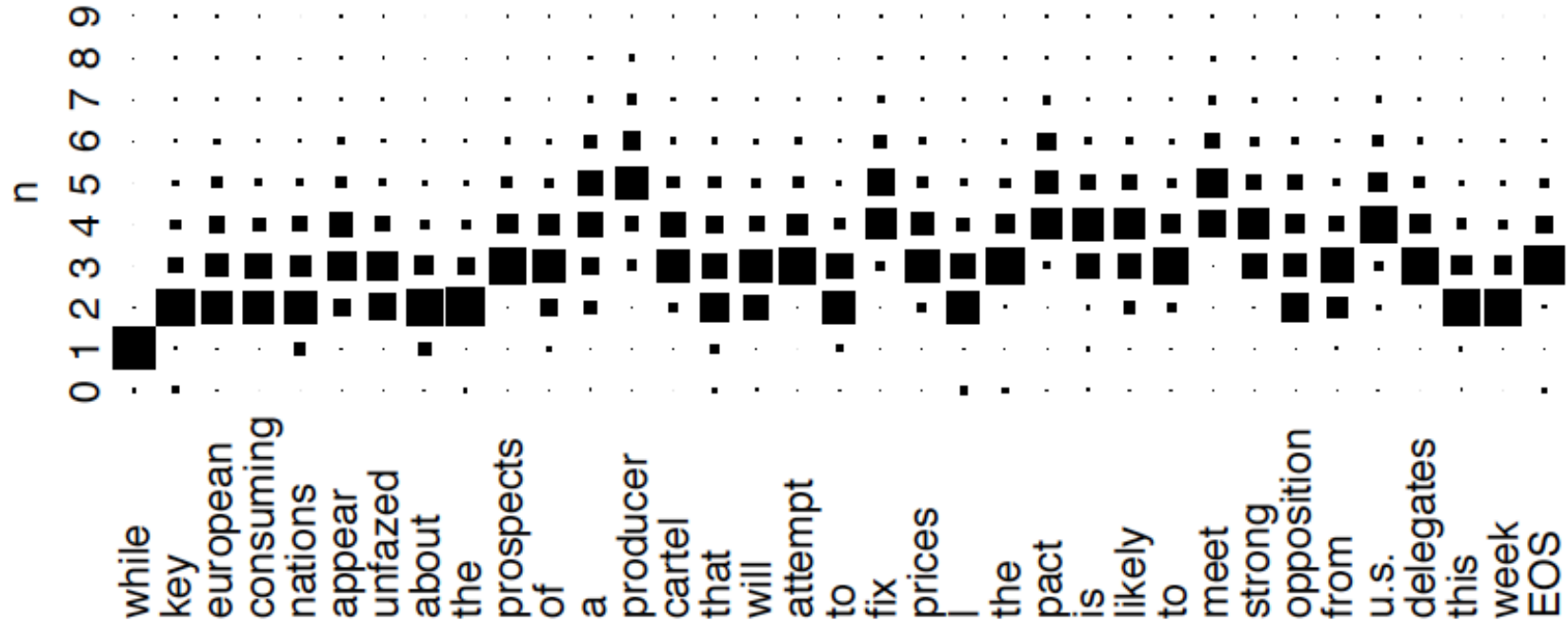
Character		s a i d _ a l i c e ; _ ‘ l e t _ u s _ a l l _ f o r _ a n y t h i n g _ t h e _ s e c o n d l y , _ . . .
Markov order		5 6 5 4 7 1 0 6 5 4 3 7 1 4 8 2 4 4 6 5 5 4 4 5 5 6 4 5 6 7 7 7 5 3 3 4 5 9 1 1 6 4 8 9 8 9 4 4 4 7 3 4 3 . . .

(b) Markov orders used to generate each character above.

Figure 5: Character-based infinite Markov model trained on “Alice in Wonderland.”

Max. order	Perplexity
$n = 3$	6.048
$n = 5$	3.803
$n = 10$	3.519
$n = \infty$	3.502

Example: Word Model



n	HPYLM	VPYLM	Nodes(H)	Nodes(V)
3	113.60	113.74	1,417K	1,344K
5	101.08	101.69	12,699K	7,466K
7	N/A	100.68	27,193K	10,182K
8	N/A	100.58	34,459K	10,434K
∞	—	100.36	—	10,629K

Modeling versus Learning

- Mochihashi & Sumita would like to learn model orders of higher order where they have more data, lower order where they have less
- They do this by making the true order finite but random under the generative model
- The Sequence Memoizer's arguably cleaner approach: Make the order infinite for all samples, but find a way to make this well regularized and computationally tractable