

The Sequence Memoizer

Frank Wood

Cedric Archambeau

Jan Gasthaus

Lancelot James

Yee Whye Teh

UCL

Gatsby

HKUST

Gatsby

Executive Summary

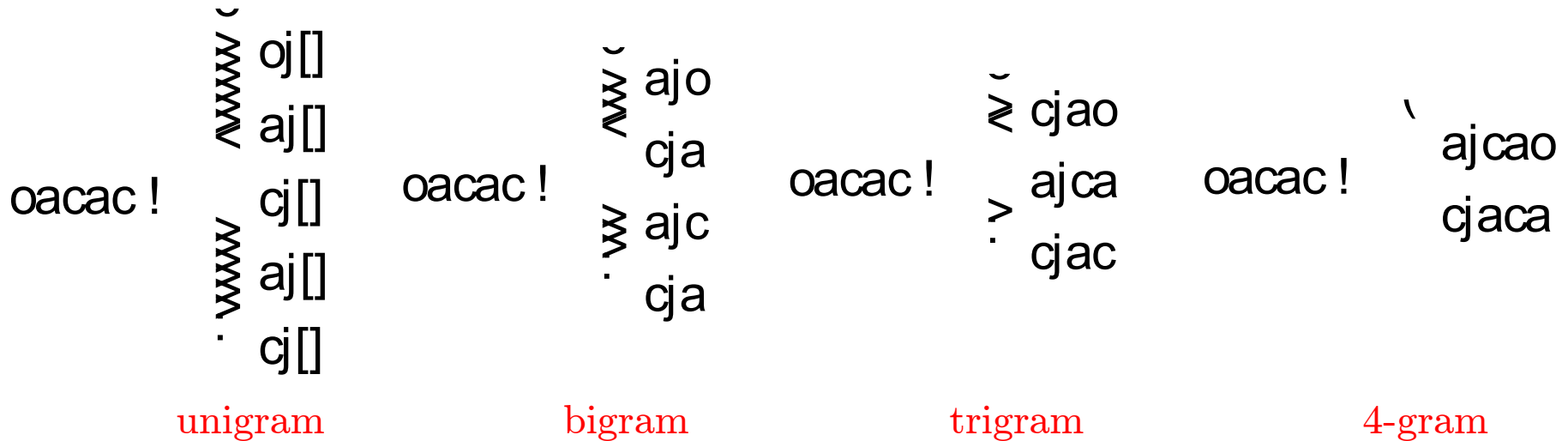
- Model
 - Smoothing Markov model of discrete sequences
 - Extension of hierarchical Pitman Yor process [Teh 2006]
 - Unbounded depth (context length)
- Algorithms and estimation
 - Linear time suffix-tree graphical model identification and construction
 - Standard Chinese restaurant franchise sampler
- Results
 - Maximum contextual information used during inference
 - Competitive language modelling results
 - Limit of n -gram language model as $n \rightarrow \infty$
 - Same computational cost as a Bayesian interpolating 5-gram language model

Executive Summary

- Uses
 - Any situation in which a low-order Markov model of discrete sequences is insufficient
 - Drop in replacement for smoothing Markov model
- Name?
 - “A Stochastic Memoizer for Sequence Data” → Sequence Memoizer (SM)
 - Describes posterior inference [Goodman et al ‘08]

Statistically Characterizing a Sequence

- Sequence Markov models are usually constructed by treating a sequence as a set of (exchangeable) observations in fixed-length contexts



Increasing context length / order of Markov model

Decreasing number of observations

Increasing number of conditional distributions to estimate (indexed by context)

Increasing power of model

Finite Order Markov Model

$$\begin{aligned}
 P(x_{1:N}) &= \prod_{i=1}^N P(x_i | x_1; \dots; x_{i-1}) \\
 &= \prod_{i=1}^N P(x_i | x_{i-n+1}; \dots; x_{i-1}); \quad n = 2 \\
 &= P(x_1)P(x_2|x_1)P(x_3|x_2)P(x_4|x_3) \dots
 \end{aligned}$$

- Example

$$\begin{aligned}
 P(oacac) &= P(o)P(a|o)P(c|a)P(a|c)P(c|a) \\
 &= G_{[o]}(o)G_{[o]}(a)G_{[c]}(a)G_{[a]}(c)G_{[c]}(a)
 \end{aligned}$$

Learning Discrete Conditional Distributions

- Discrete distribution \leftrightarrow vector of parameters

$$G_{[u]} = [1/4; \dots; 1/k]; K \geq j \geq 1$$

- Counting / Maximum likelihood estimation

- Training sequence $x_{1:N}$

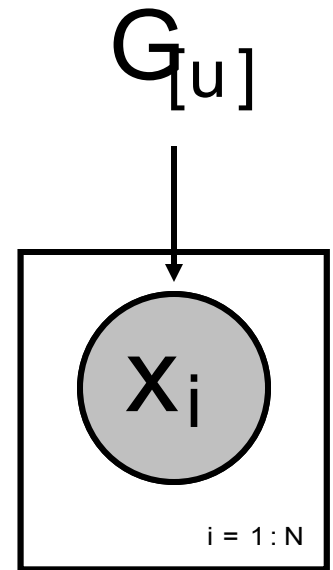
$$\hat{G}_{[u]}(X = k) = \hat{\pi}_k = \frac{\# \text{f uk g}}{\# \text{f u g}}$$

- Predictive inference

$$P(X_{n+1} | x_1 : : : x_N) = \hat{G}_{[u]}(X_{n+1})$$

- Example

- Non-smoothed unigram model ($\mathbf{u} = \epsilon$)



Bayesian Smoothing

- Estimation

$$P(G_{[u]} | x_{1:n}) / P(x_{1:n} | G_{[u]}) P(G_{[u]})$$

- Predictive inference

$$P(X_{n+1} | x_{1:n}) = \int P(X_{n+1} | G_{[u]}) P(G_{[u]} | x_{1:n}) dG_{[u]}$$

- Priors over distributions

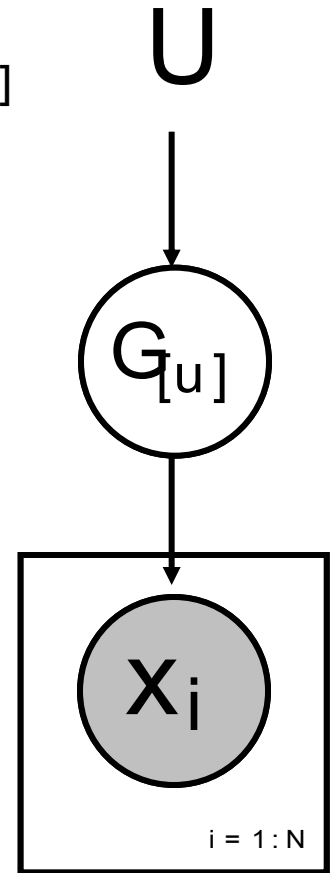
$$G_{[u]} \gg \text{Dirichlet}(U); \quad G_{[u]} \gg \text{PY}(d; c; U)$$

- Net effect

- Inference is “smoothed” w.r.t. uncertainty about unknown *distribution*

- Example

- Smoothed unigram ($\mathbf{u} = \epsilon$)



A Way To Tie Together Distributions

$$\begin{array}{l}
 G_{[u]} \quad \gg \quad \text{PY} (d; c; G_{[\frac{3}{4}(u)]}) \\
 X_i \quad \gg \quad G_{[u]} \quad \text{base distribution}
 \end{array}$$

discount concentration

- Tool for tying together related distributions in hierarchical models
- Measure over measures
- Base measure is the “mean” measure


$$E[G_{[u]}(dx)] = G_{[\frac{3}{4}(u)]}(dx)$$

- A distribution drawn from a Pitman Yor process is related to its base distribution
 - (equal when $c = \infty$ or $d = 1$)

Pitman-Yor Process Continued

- Generalization of the Dirichlet process ($d = 0$)
 - Different (power-law) properties
 - Better for text [Teh, 2006] and images [Sudderth and Jordan, 2009]
- Posterior predictive distribution

$$\begin{aligned}
 P(X_{N+1} | X_{1:N}; c; d) &= \int P(X_{N+1} | G_u) P(G_u | X_{1:N}; c; d) dG_u \\
 &= E \left[\frac{\sum_{k=1}^K (m_k + d) \mathbb{1}(\hat{A}_k = X_{N+1})}{c + N} + \frac{c + dK}{c + N} G_{[3/4(u)]}(X_{N+1}) \right]
 \end{aligned}$$

Can't actually do this integral this way 

- Forms the basis for straightforward, simple samplers
- Rule for stochastic memoization

Hierarchical Bayesian Smoothing

- Estimation

$$\mathcal{E} = f(G_{[u]}; G_{[v]}; G_{[w]})g; \quad w = \frac{3}{4}(u) = \frac{3}{4}(v)$$

$$P(\mathcal{E} | x_{1:N}) / P(x_{1:N} | \mathcal{E}) P(\mathcal{E})$$

- Predictive inference

$$P(X_{N+1} | x_{1:N})$$

$$= \int P(X_{N+1} | \mathcal{E}) P(\mathcal{E} | x_{1:N}) d\mathcal{E}$$

- Naturally related distributions tied together

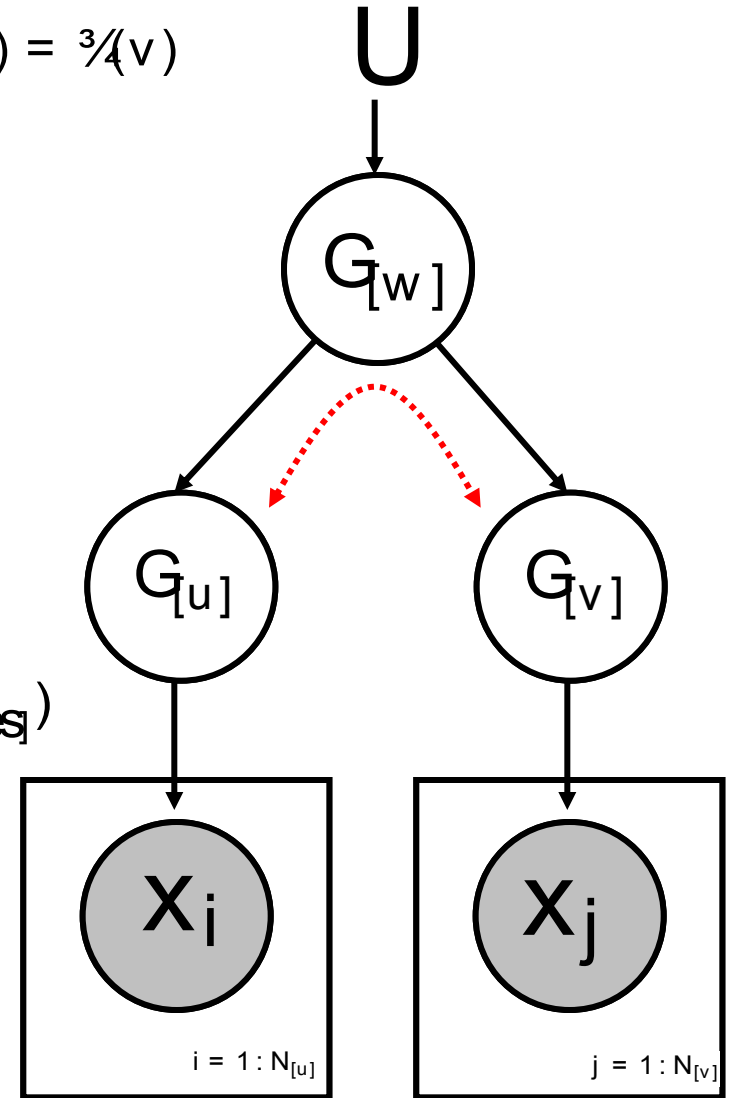
$$G_{[\text{the United States}]} \gg PY(d; c; G_{[\text{United States}]})$$

- Net effect

- Observations in one context affect inference in other context.
- Statistical strength is shared between similar contexts

- Example

- Smoothing bi-gram ($w = \epsilon, u, v \in \Sigma$)

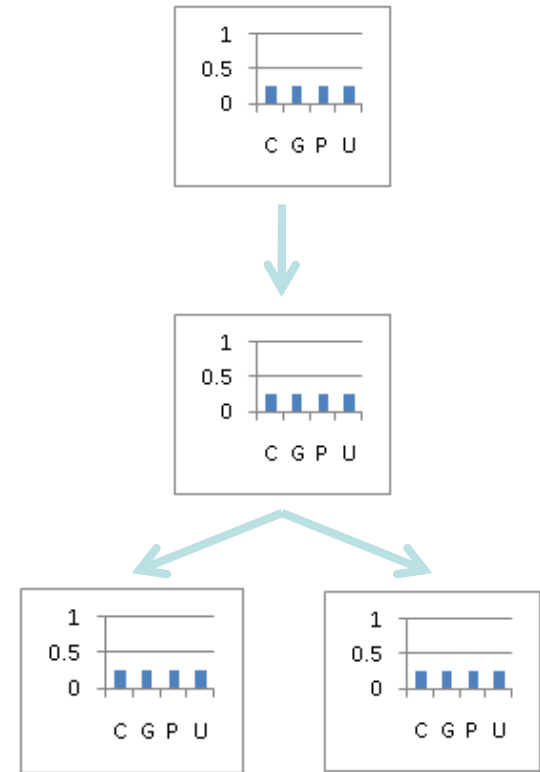
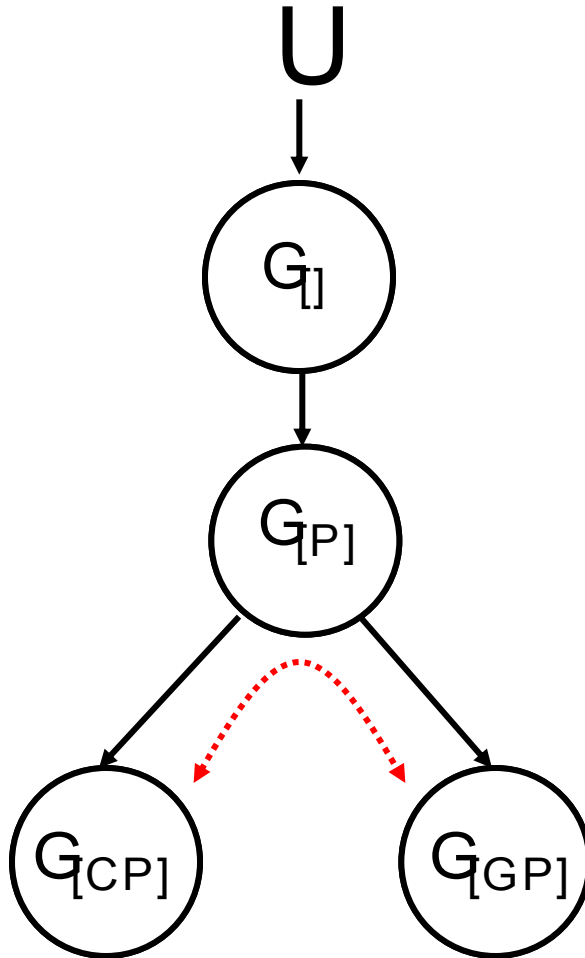


SM/HPYP Sharing in Action

Observations

Conditional Distributions

Posterior Predictive Probabilities



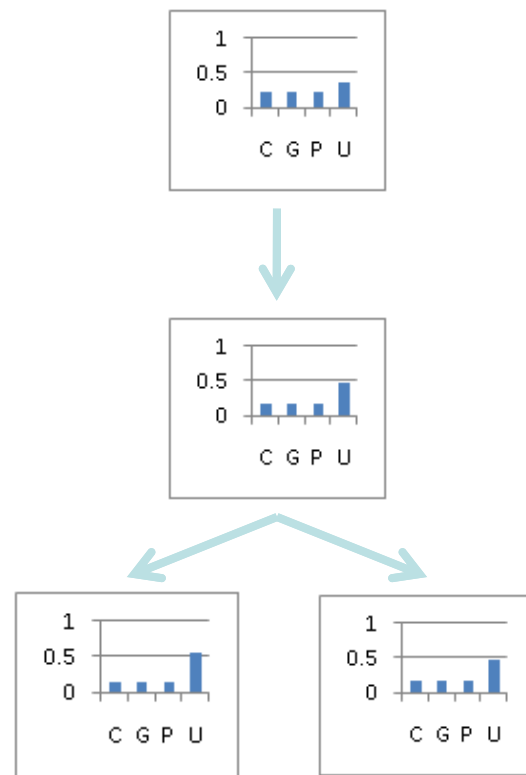
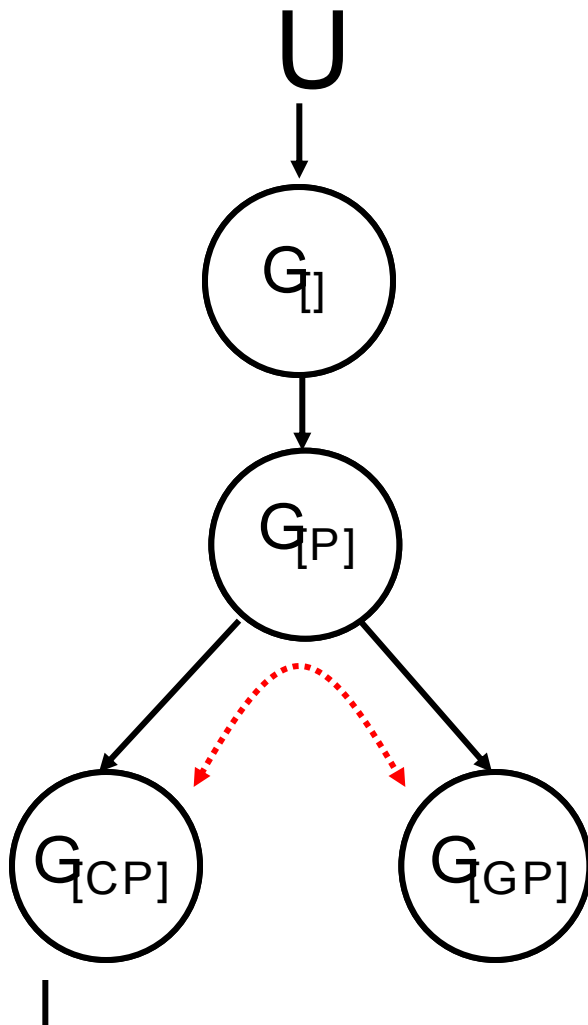
CRF Particle Filter Posterior Update

Observations

Conditional Distributions

Posterior Predictive Probabilities

CPU



CRF Particle Filter Posterior Update

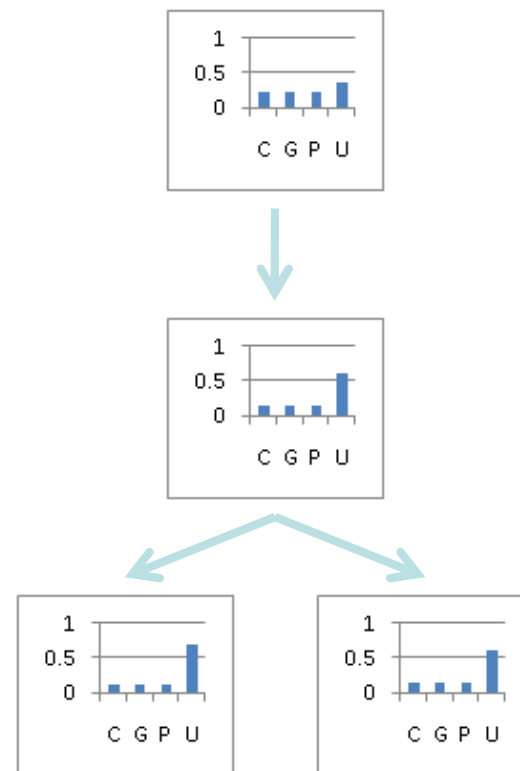
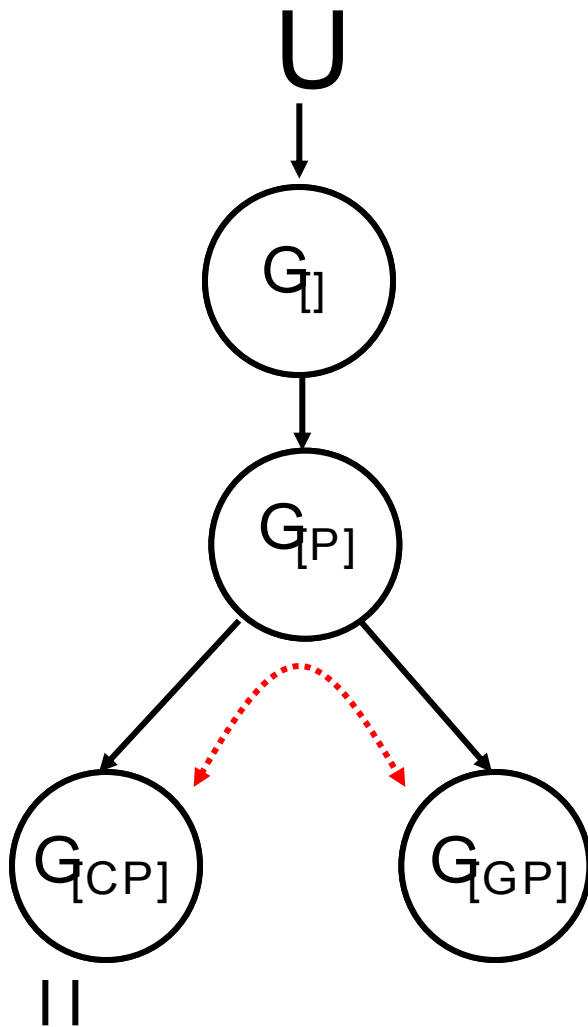
Observations

Conditional Distributions

Posterior Predictive Probabilities

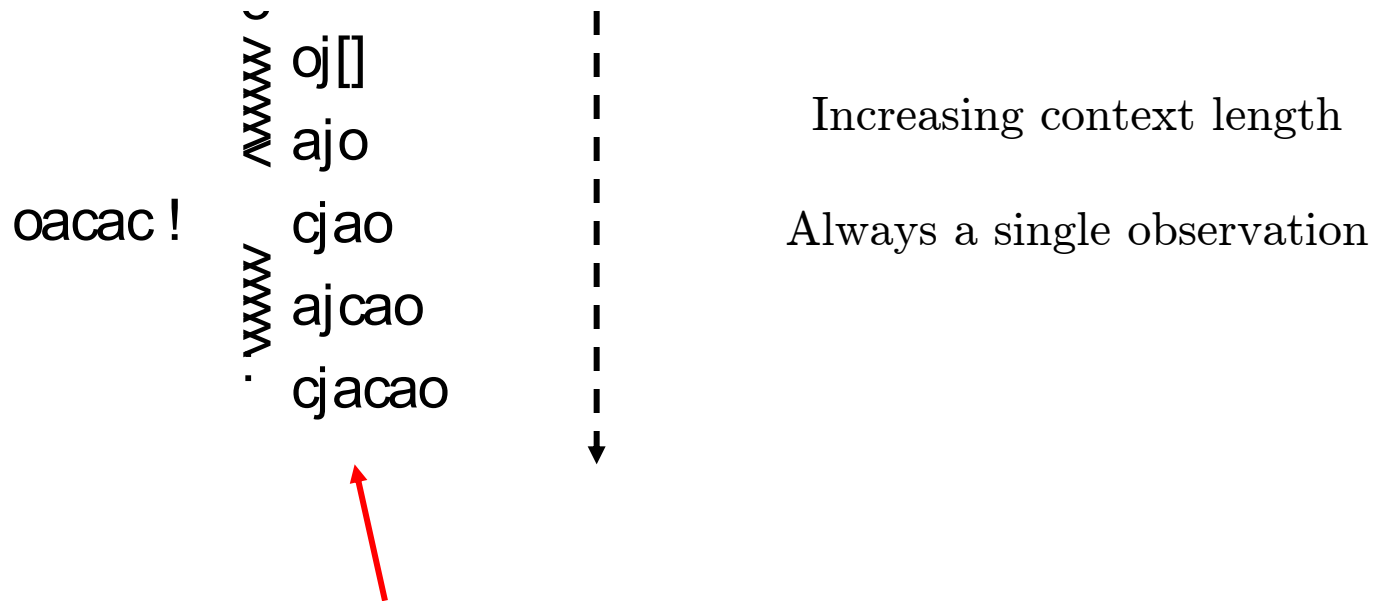
CPU

CPU



Alternative Sequence Characterization

- A sequence can be characterized by a set of *single* observations in unique contexts of growing length



Foreshadowing: all suffixes of the string "cacao"

“Non-Markov” Model

$$\begin{aligned} P(x_{1:N}) &= \prod_{i=1}^N P(x_i | x_1; \dots; x_{i-1}) \\ &= P(x_1)P(x_2|x_1)P(x_3|x_2; x_1)P(x_4|x_3; \dots; x_1) \dots \end{aligned}$$

- Example

$$P(oacac) = P(o)P(a|o)P(c|oa)P(a|oac)P(c|oaca)$$

- Smoothing essential

- Only one observation in each context!

- Solution

- Hierarchical sharing ala HPYP

HPYP LM Sharing Architecture

- Share statistical strength between sequentially related predictive conditional distributions

- Estimates of highly specific conditional distributions

$G_{[was\ on\ the]}$

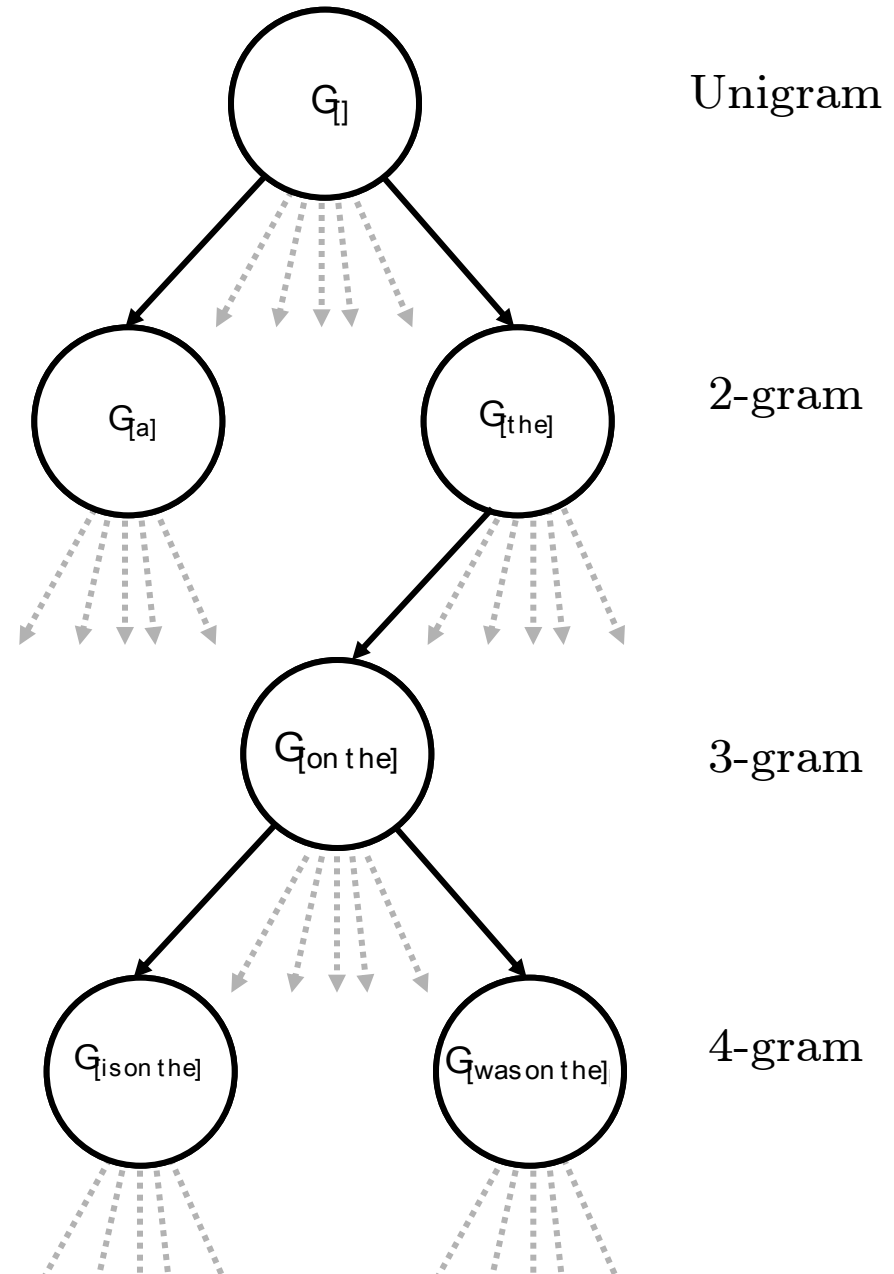
- Are coupled with others that are related

$G_{[is\ on\ the]}$

- Through a single common, more-general shared ancestor

$G_{[on\ the]}$

- Corresponds intuitively to back-off



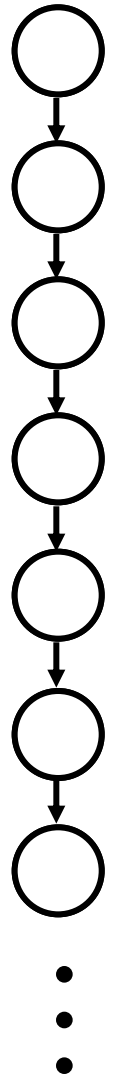
Hierarchical Pitman Yor Process

$$\begin{aligned}
 G_{[j]} & \mid d_0; U && \gg && \text{PY}(d_0; 0; U) \\
 G_{[u]} & \mid d_{j_{u|j}}; G_{[\frac{3}{4}(u)]} && \gg && \text{PY}(d_{j_{u|j}}; 0; G_{[\frac{3}{4}(u)]}) \\
 x_i & \mid x_{1:i-1} = u && \gg && G_{[u]} \\
 i & = 1; \dots; T \\
 \delta_u & \propto \xi^{n_i - 1}
 \end{aligned}$$

- Bayesian generalization of smoothing n -gram Markov model
- Language model : outperforms interpolated Kneser-Ney (KN) smoothing
- Efficient inference algorithms exist
 - [Teh, '06; Teh, Kurihara, Welling, '08]
- Sharing between contexts that differ in most distant symbol only
- Finite depth

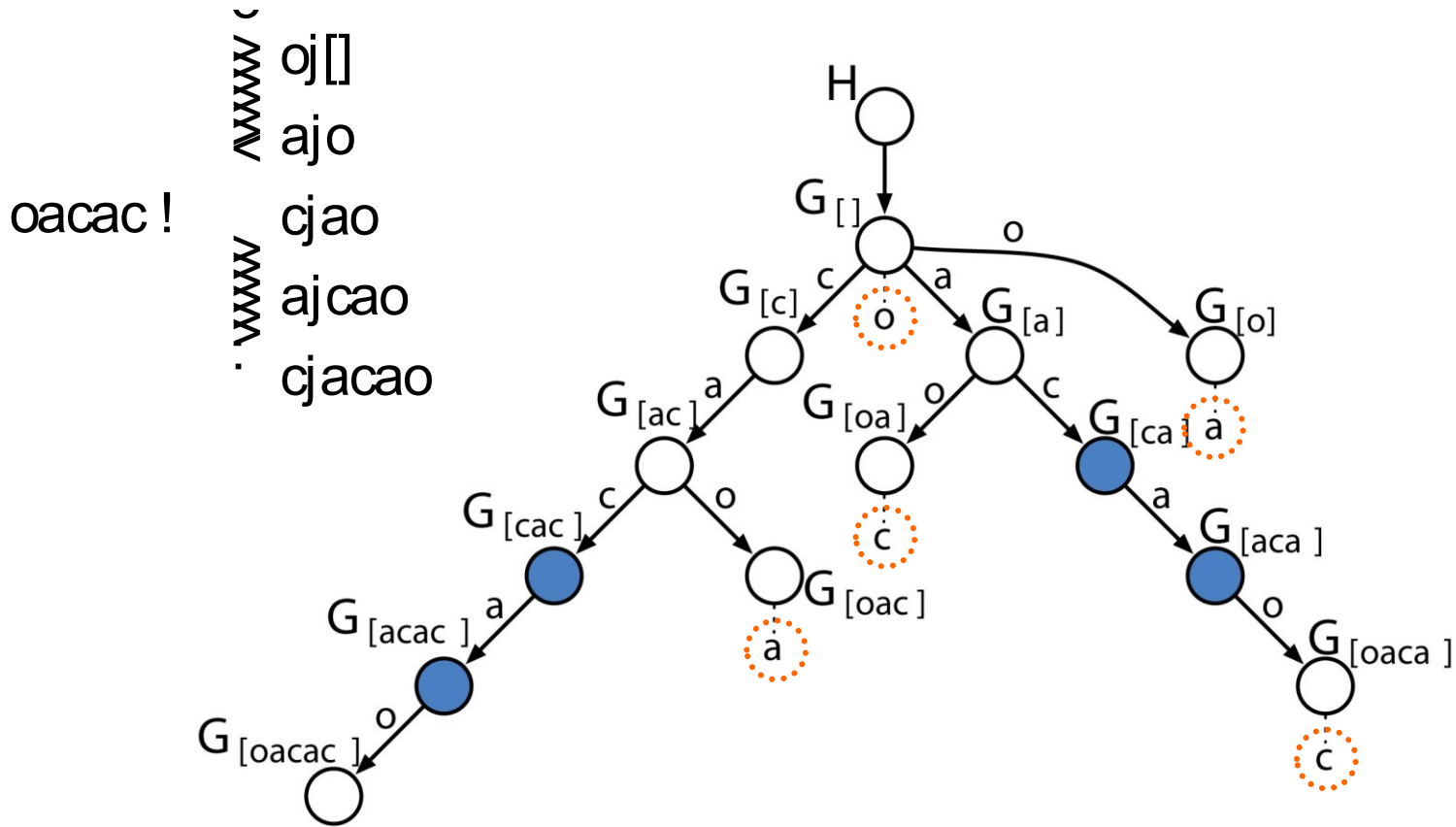
Sequence Memoizer

$$\begin{aligned}
 G_{[j]} & \text{ j } d_0; U & \gg & \text{PY}(d_0; 0; U) \\
 G_{[u]} & \text{ j } d_{j|u}; G_{[3/4(u)]} & \gg & \text{PY}(d_{j|u}; 0; G_{[3/4(u)]}) \\
 x_i & \text{ j } x_{1:i}; 1 = u & \gg & G_{[u]} \\
 & i = 1; \dots; T & & \\
 & & & \text{8u 2 } \xi^+
 \end{aligned}$$



- Eliminates Markov order selection
- Always uses full context when making predictions
- Linear time, linear space (in length of observation sequence) graphical model identification
- Performance is limit of n -gram as $n \rightarrow \infty$
- Same or less overall cost as 5-gram interpolating Kneser Ney

Graphical Model Trie



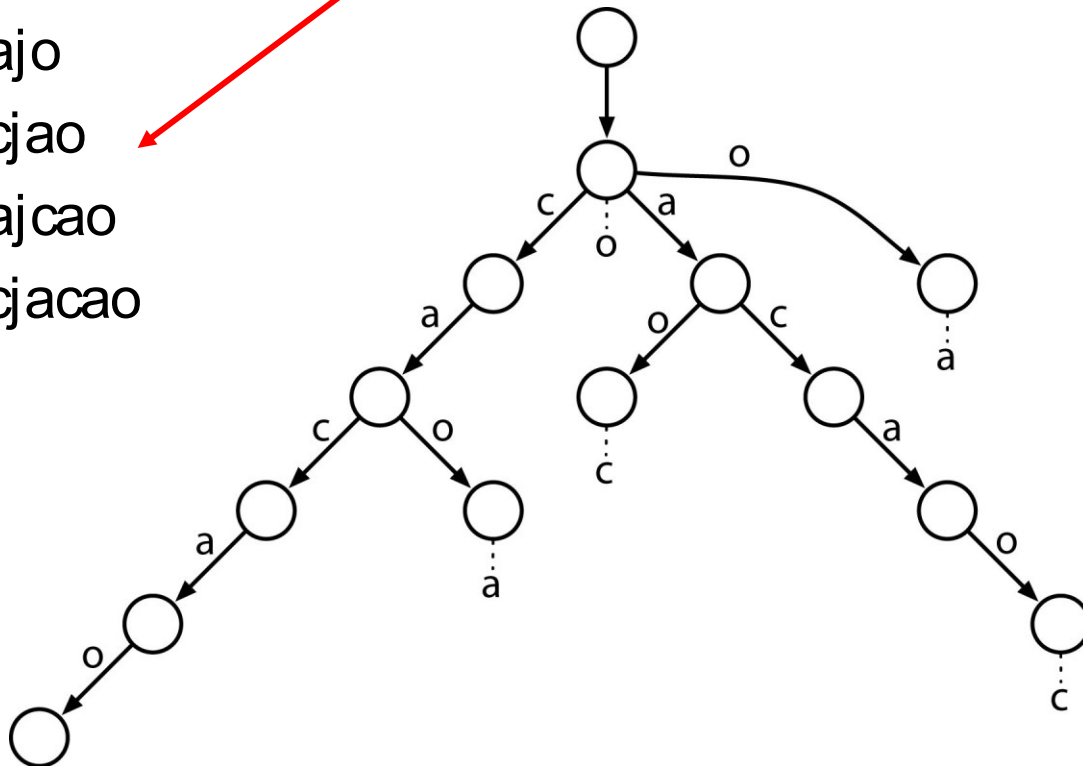
-  Observations
-  Latent conditional distributions with Pitman Yor priors / stochastic memoizers

Suffix Trie Datastructure

All suffixes of the string "cacao"

oacac!

- o j []
- a j o
- c j a o
- a j c a o
- c j a c a o



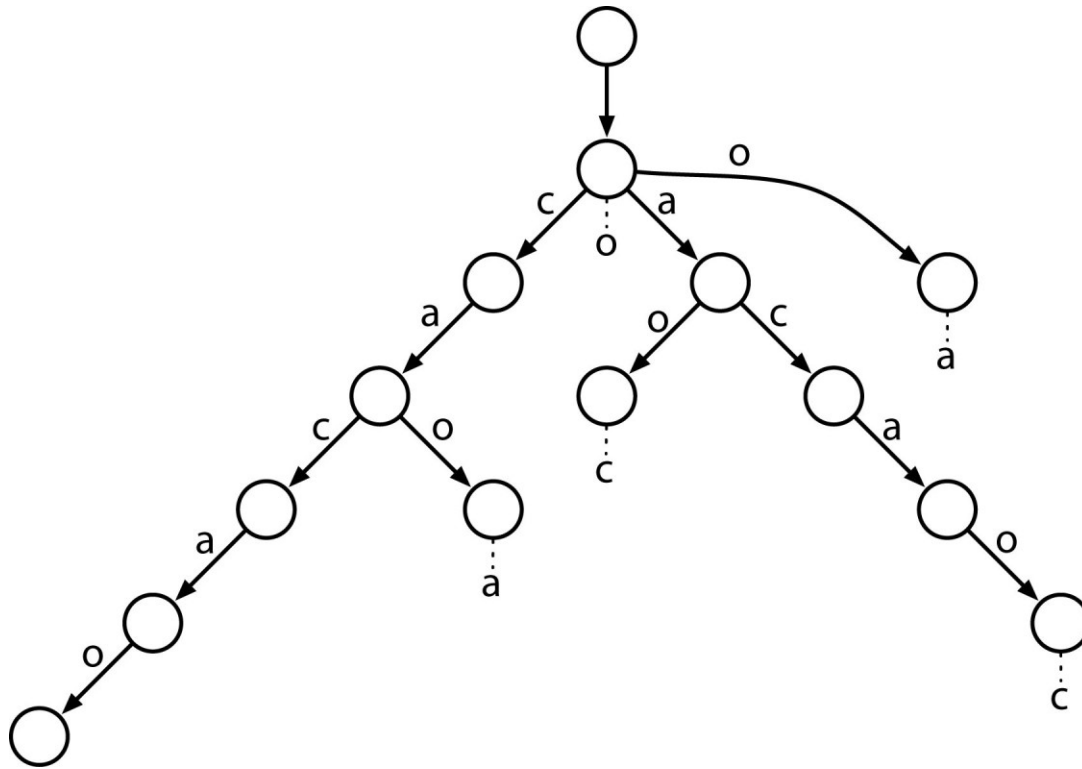
Suffix Trie Datastructure

- Deterministic finite automata that recognizes all suffixes of an input string.
- Requires $O(N^2)$ time and space to build and store [Ukkonen, 95]
- Too intensive for any practical sequence modelling application.

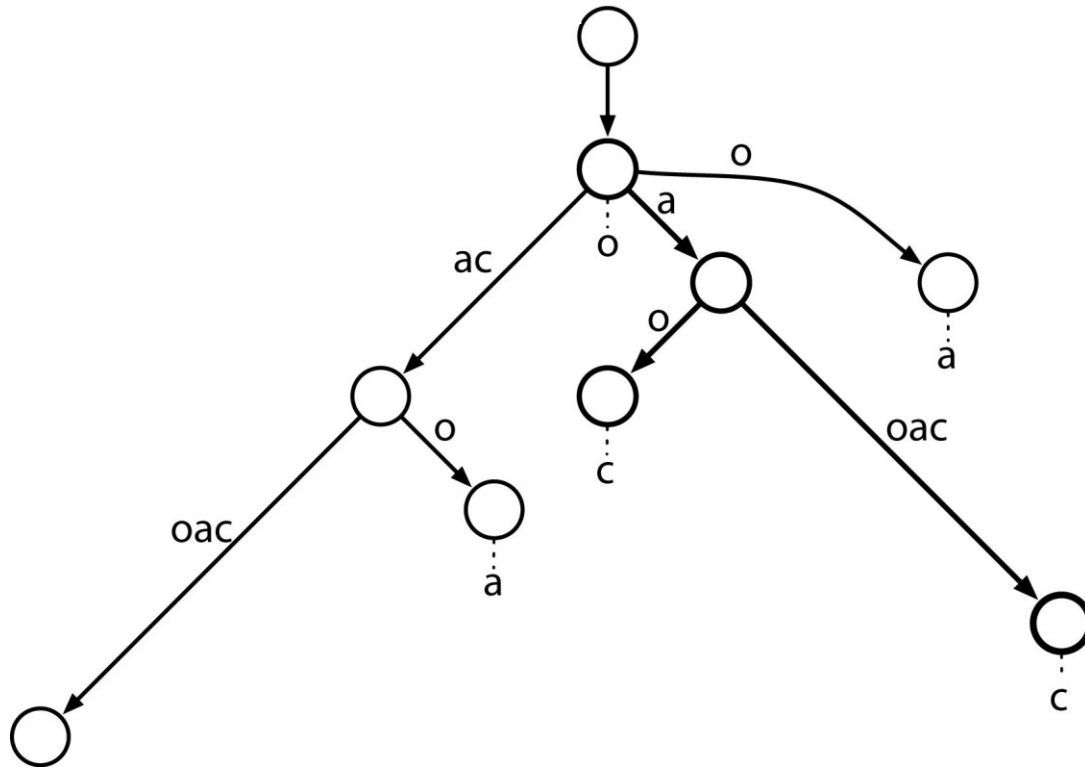
Suffix Tree

- Deterministic finite automata that recognizes all suffixes of an input string
- Uses path compression to reduce storage and construction computational complexity.
- Requires only $O(N)$ time and space to build and store [Ukkonen, 95]
- Practical for large scale sequence modelling applications

Suffix Trie Datastructure



Suffix Tree Datastructure



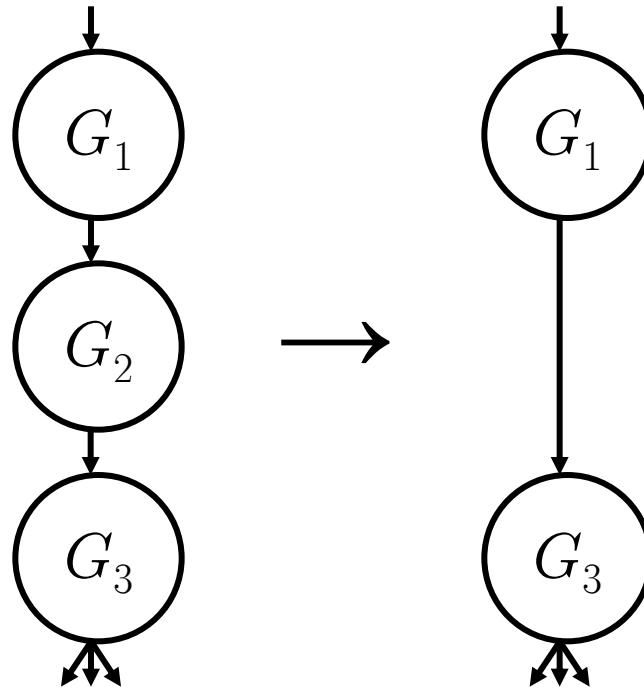
Graphical Model Identification

- This is a graphical model transformation under the covers.
- These compressed paths require being able to analytically marginalize out nodes from the graphical model
- The result of this marginalization can be thought of as providing a different set of caching rules to memoizers on the path-compressed edges

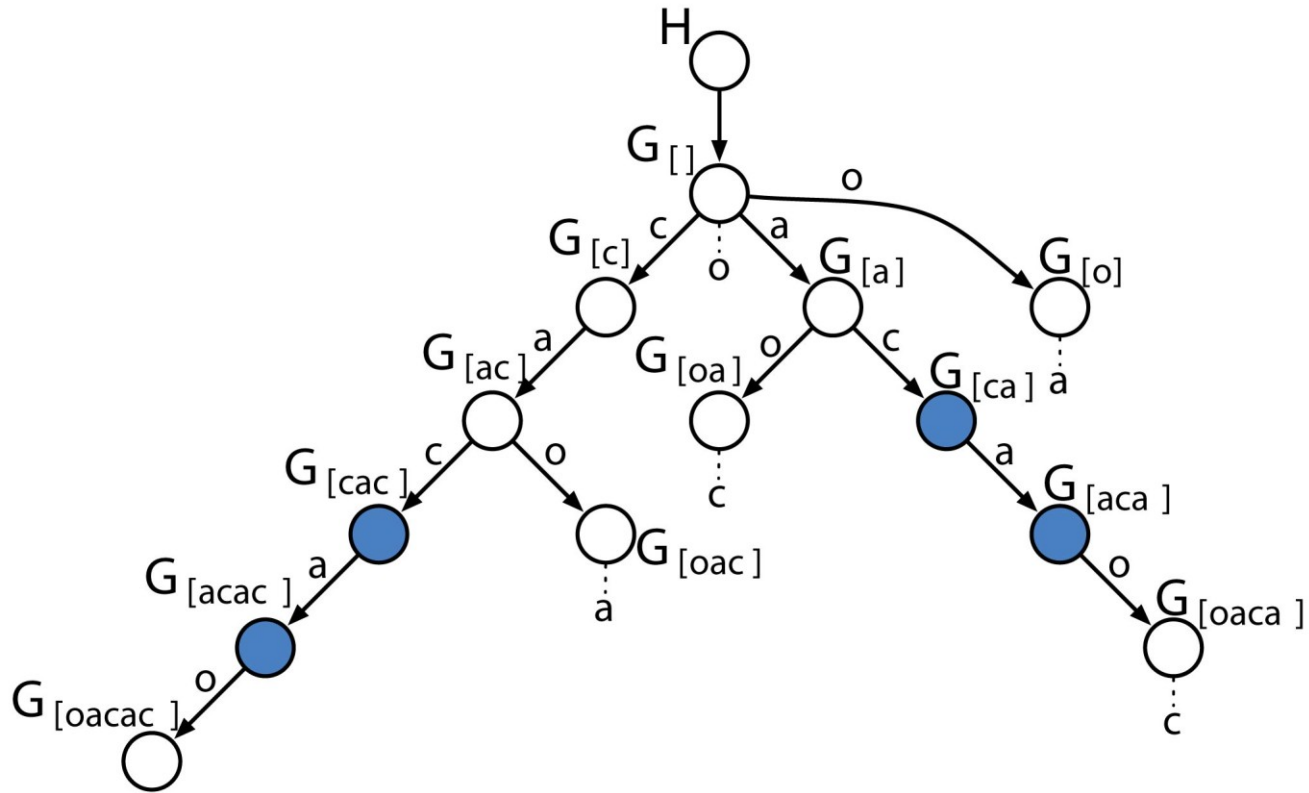
Marginalization

- Theorem 1: Coagulation

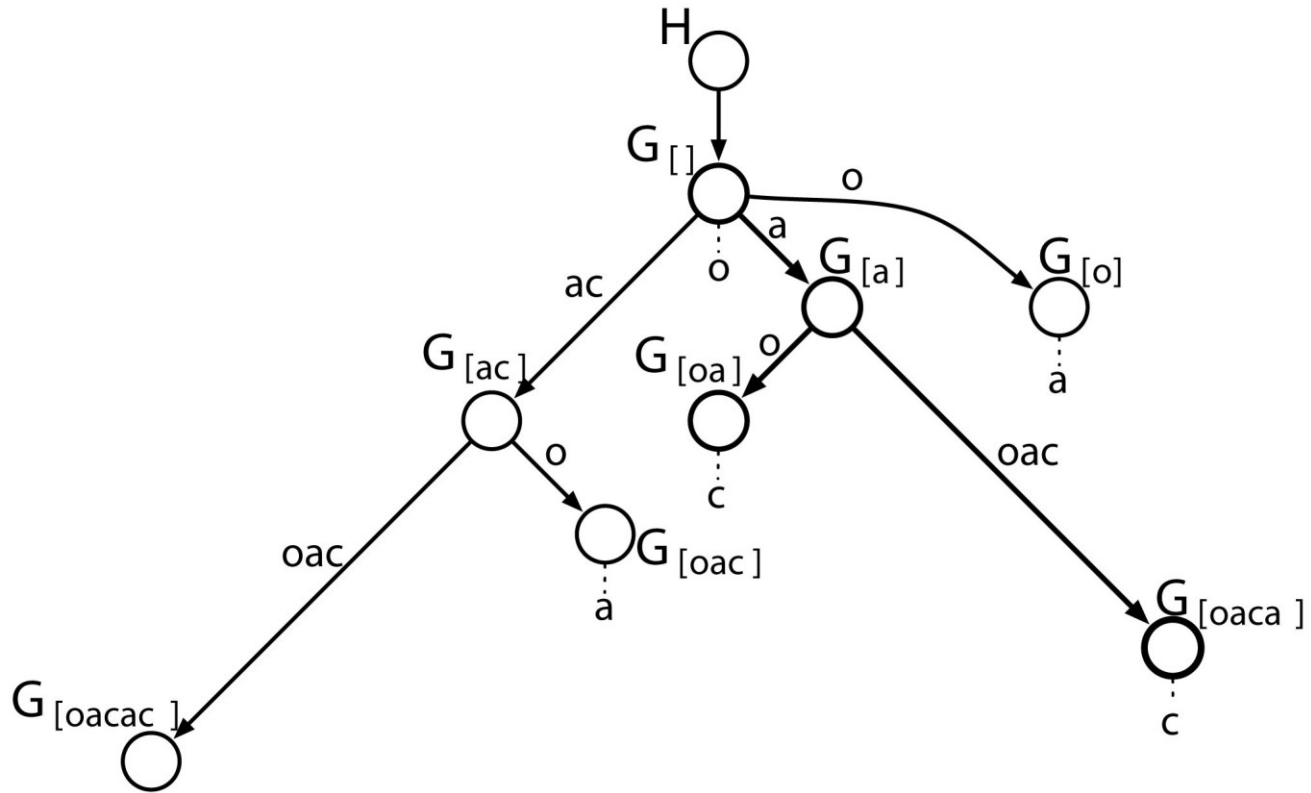
If $G_2 \downarrow G_1 \gg \text{PY}(d_1; 0; G_1)$ and $G_3 \downarrow G_2 \gg \text{PY}(d_2; 0; G_2)$
then $G_3 \downarrow G_1 \gg \text{PY}(d_1 d_2; 0; G_1)$ with G_2 marginalized out.



Graphical Model Trie



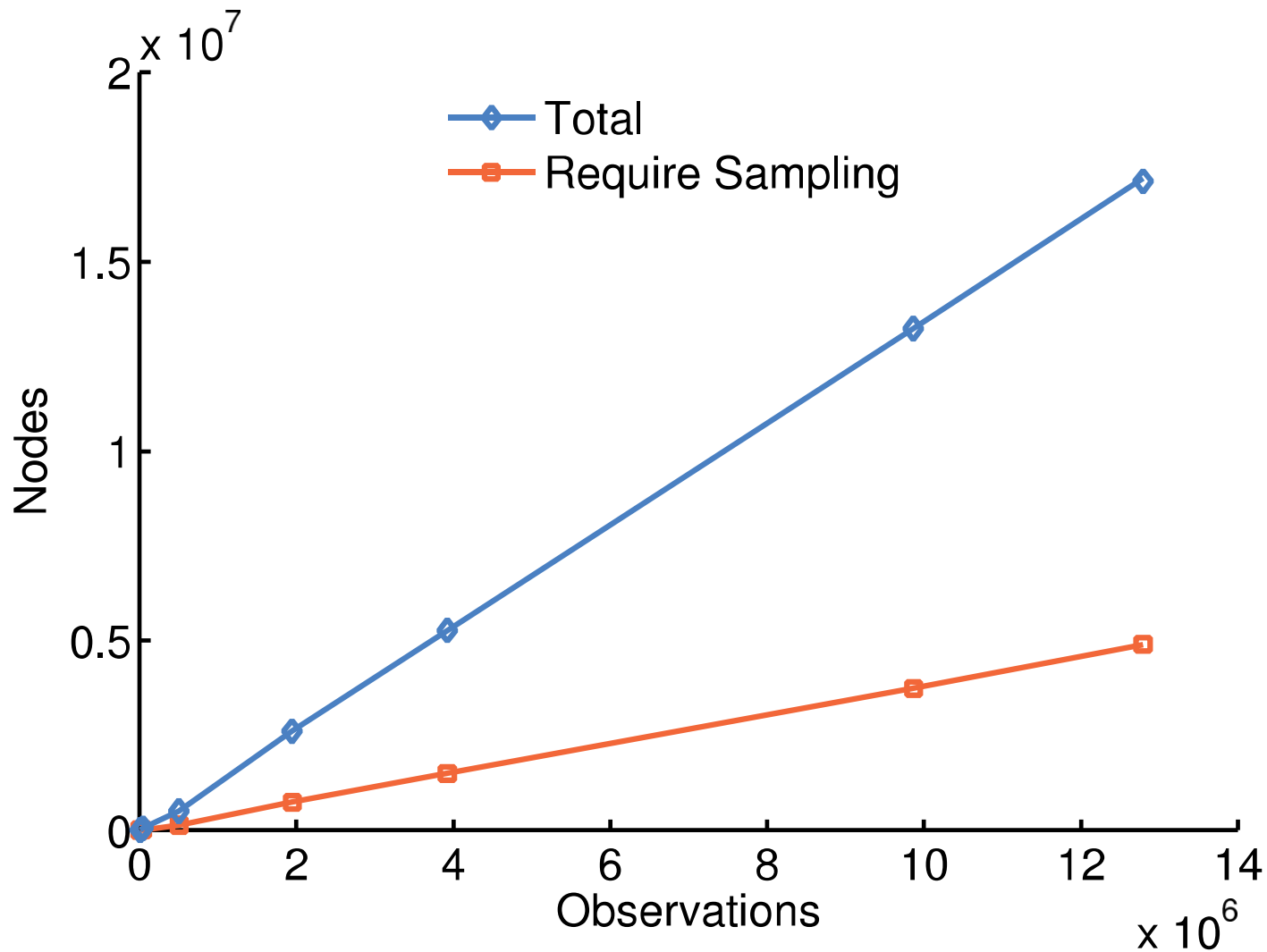
Graphical Model Tree



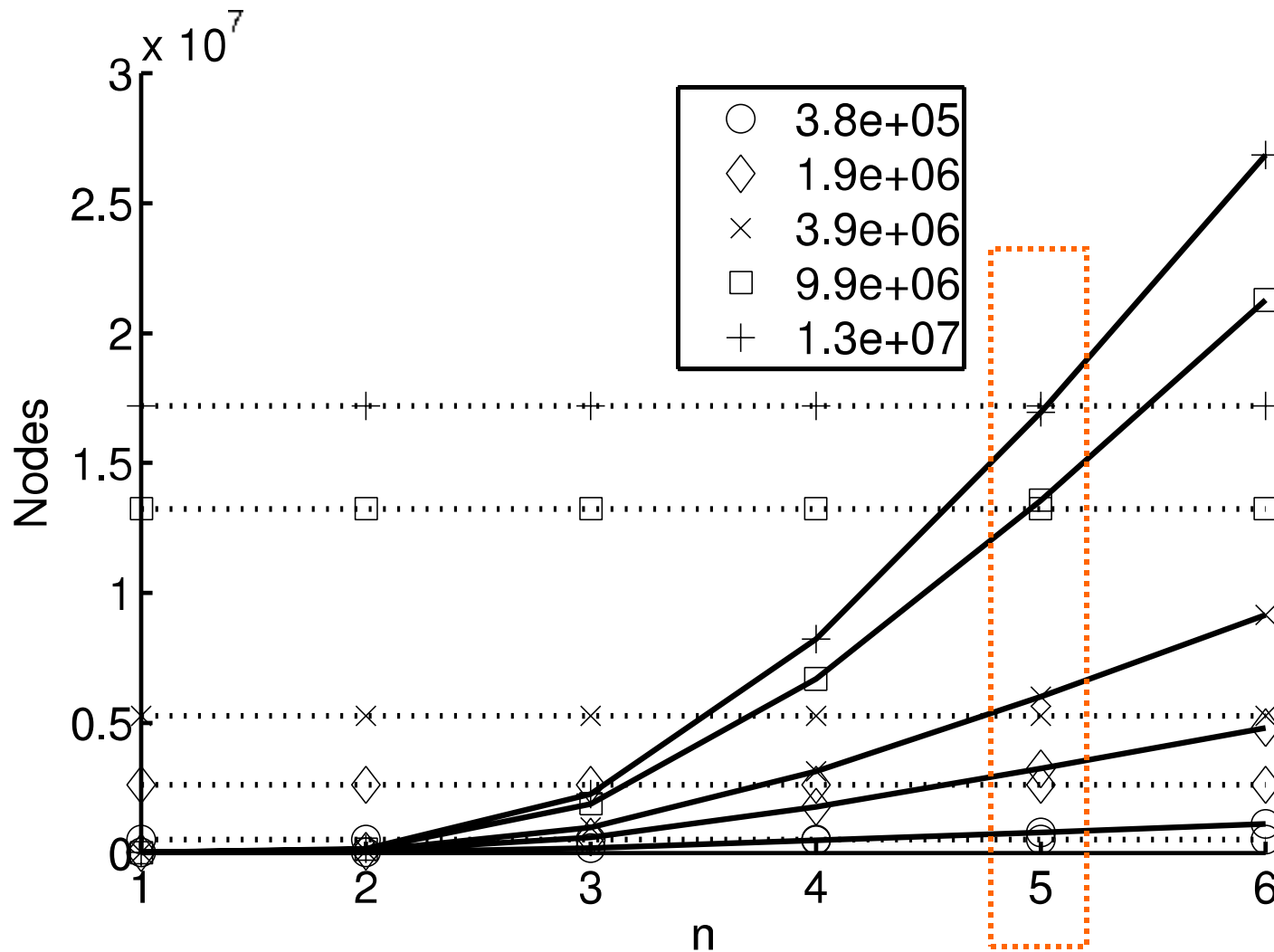
Graphical Model Initialization

- Given a single input sequence
 - Ukkonen's linear time suffix tree construction algorithm is run on its reverse to produce a prefix tree
 - This identifies the nodes in the graphical model we need to represent
 - The tree is traversed and path compressed parameters for the Pitman Yor processes are assigned to each remaining Pitman Yor process

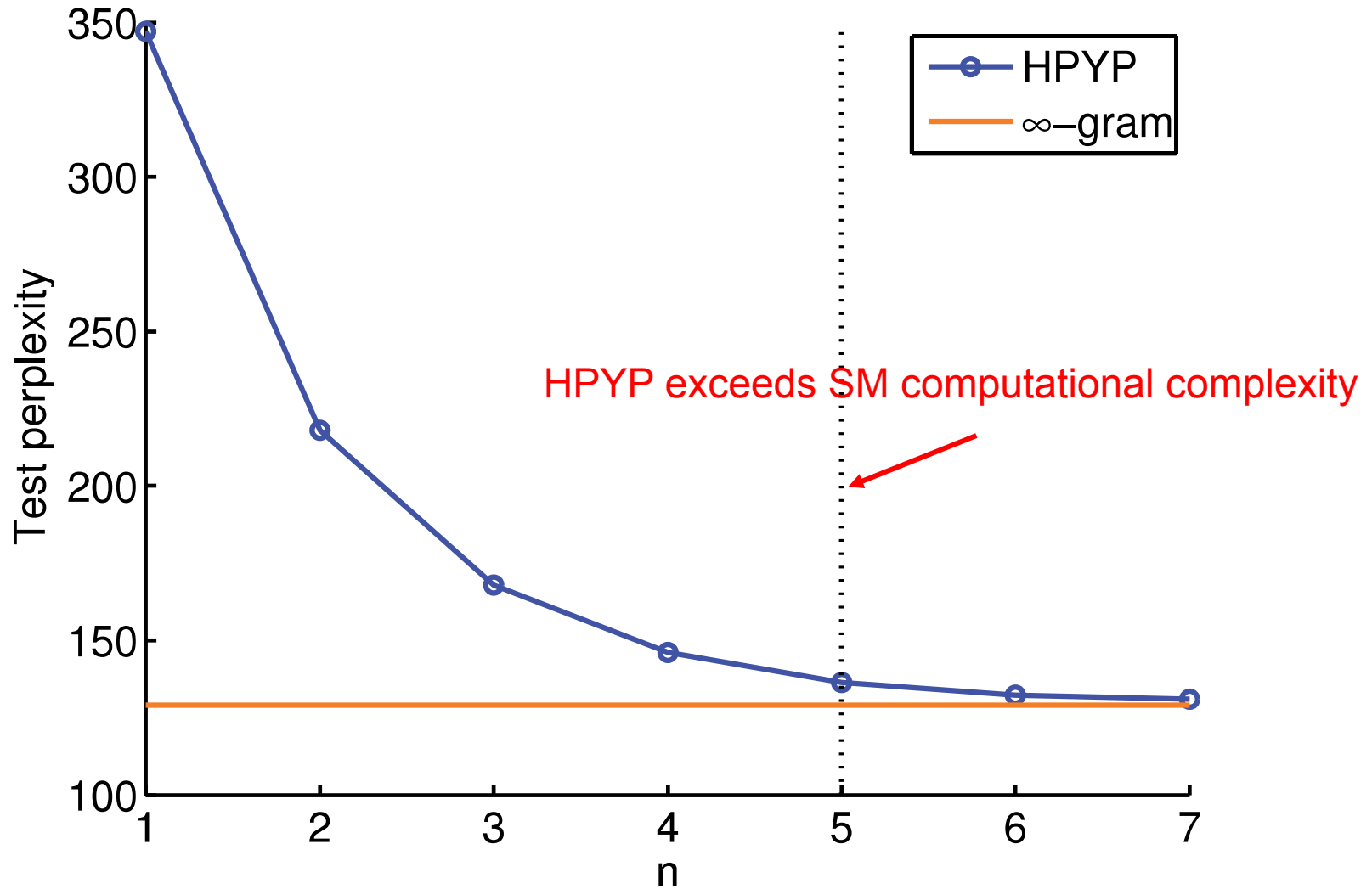
Nodes In The Graphical Model



Never build more than a 5-gram



Sequence Memoizer Bounds N-Gram Performance



Language Modelling Results

AP News Test Perplexity

[Mnih & Hinton, 2009]	112.1
[Bengio et al., 2003]	109.0
4-gram Modified Kneser-Ney [Teh, 2006]	102.4
4-gram HPYP [Teh, 2006]	101.9
Sequence Memoizer (SM)	96.9

The Sequence Memoizer

- The Sequence Memoizer is a deep (unbounded) smoothing Markov model
- It can be used to learn a joint distribution over discrete sequences in time and space linear in the length of a single observation sequence
- It is equivalent to a smoothing ∞ -gram but costs no more to compute than a 5-gram