

# Improving nonparameteric Bayesian inference: experiments on unsupervised word segmentation with adaptor grammars

Mark Johnson<sup>1</sup>   Sharon Goldwater<sup>2</sup>

<sup>1</sup>Brown University  
Providence, RI  
Mark\_Johnson@Brown.edu

<sup>2</sup>University of Edinburgh  
Edinburgh EH8 9AB  
sgwater@inf.ed.ac.uk

NAACL 2009 Conference

# Outline

- *Non-parametric models* learn which parameters are relevant
  - ▶ Unsupervised word segmentation: identify the words in an unsegmented stream of phonemes
- *Adaptor grammars* are a non-parametric extension of PCFGs which learn the rules as well as their probabilities
- Monte Carlo Markov Chain (MCMC) *sampling methods* are natural inference procedures for non-parametric models
  - ▶ The adaptor grammar inference procedure samples rules as well their probabilities
- This paper describes improvements to adaptor grammars and their MCMC inference procedure
  - ▶ *estimating hyper-parameters* using slice sampling
  - ▶ *modal decoding* from multiple samples from multiple runs
  - ▶ *random initialization* instead of incremental initialization
  - ▶ *table label resampling* as well as sentence resampling

All but the last apply to a wide variety of models

# Outline

Non-parametric inference in word segmentation

Chinese Restaurant Processes

Adaptor grammars

- Adaptor grammars for word segmentation

- Priors on Adaptor Grammar PYP parameters

Unsupervised inference of Adaptor Grammars

- Modal word segmentation

- Random vs incremental initialization

- Table label resampling

Conclusion

# Unsupervised word segmentation

- Input: phoneme sequences with *sentence boundaries* (Brent)
- Task: identify *word boundaries*, and hence words

y Δ u Δ w Δ a Δ n Δ t Δ t Δ u Δ s Δ i Δ D Δ 6 Δ b Δ U Δ k

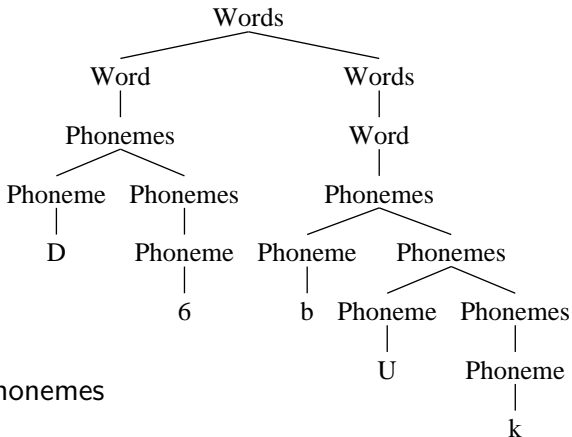
- Useful cues for word segmentation:
  - ▶ Phonotactics (Fleck)
  - ▶ Inter-word dependencies (Goldwater)

# Word segmentation with PCFGs (1)

Sentence  $\rightarrow$  Word<sup>+</sup>  
Word  $\rightarrow$  Phoneme<sup>+</sup>

which abbreviates

Sentence  $\rightarrow$  Words  
Words  $\rightarrow$  Word Words  
Word  $\rightarrow$  Phonemes  
Phonemes  $\rightarrow$  Phoneme Phonemes  
Phonemes  $\rightarrow$  Phoneme  
Phoneme  $\rightarrow a \mid \dots \mid z$

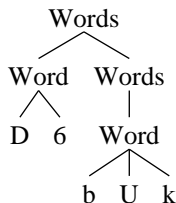


# Word segmentation with PCFGs (1)

Sentence  $\rightarrow$  Word<sup>+</sup>

Word  $\rightarrow$  all possible phoneme strings

- But now there are an infinite number of PCFG rules!
  - ▶ once we see our (finite) training data, only finitely many are useful
  - $\Rightarrow$  the set of parameters (rules) should be chosen based on training data



# Outline

Non-parametric inference in word segmentation

## Chinese Restaurant Processes

Adaptor grammars

- Adaptor grammars for word segmentation

- Priors on Adaptor Grammar PYP parameters

Unsupervised inference of Adaptor Grammars

- Modal word segmentation

- Random vs incremental initialization

- Table label resampling

Conclusion

# Multinomials and Chinese Restaurant Processes

- Multinomials (e.g., dice, HMMs, PCFGs)
  - ▶ Observations  $\mathbf{z} = (z_1, \dots, z_n)$  ranging over outcomes  $1, \dots, m$  (e.g., sides of a die), with outcome  $k$  observed  $n_k(\mathbf{z})$  times
  - ▶ *Predictive distribution* with uniform Dirichlet prior:

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto n_k(\mathbf{z}) + \alpha/m$$

- Let  $m \rightarrow \infty$

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto n_k(\mathbf{z}) \text{ if } k \text{ appears in } \mathbf{z}$$

$$P(Z_{n+1} \notin \mathbf{z} \mid \mathbf{z}) \propto \alpha$$

- Insist outcomes are enumerated *in order*  
 $\Rightarrow$  *Chinese Restaurant Process*

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) & \text{if } k \leq m = \max(\mathbf{z}) \\ \alpha & \text{if } k = m + 1 \end{cases}$$

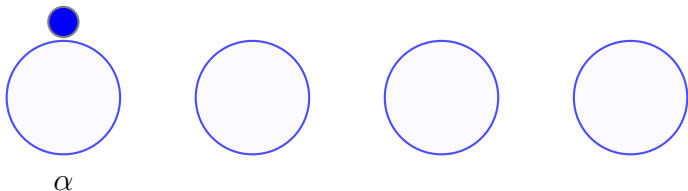


# Chinese Restaurant Process (0)



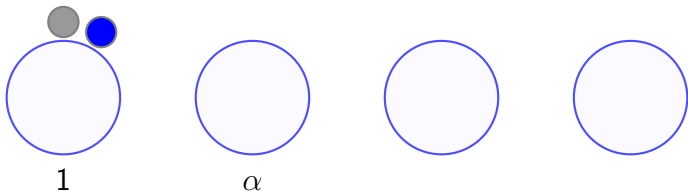
- Customer  $\rightarrow$  table mapping  $\mathbf{z} =$
- $P(\mathbf{z}) = 1$
  
- Next customer chooses a table according to:
  - $P(\text{occupied table}) \propto \text{number of customers at table}$
  - $P(\text{next unoccupied table}) \propto \alpha$

# Chinese Restaurant Process (1)



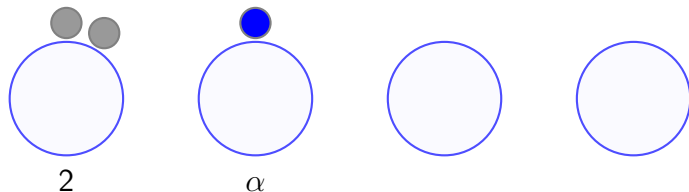
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1$
- $P(\mathbf{z}) = \alpha/\alpha$
- Next customer chooses a table according to:
  - $P(\text{occupied table}) \propto \text{number of customers at table}$
  - $P(\text{next unoccupied table}) \propto \alpha$

## Chinese Restaurant Process (2)



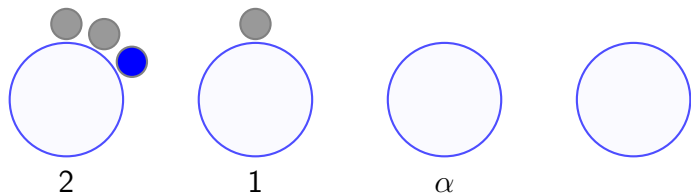
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1$
- $P(\mathbf{z}) = \alpha/\alpha \times 1/(1 + \alpha)$
- Next customer chooses a table according to:
  - $P(\text{occupied table}) \propto \text{number of customers at table}$
  - $P(\text{next unoccupied table}) \propto \alpha$

## Chinese Restaurant Process (3)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2$
- $P(\mathbf{z}) = \alpha/\alpha \times 1/(1 + \alpha) \times \alpha/(2 + \alpha)$
- Next customer chooses a table according to:
  - $P(\text{occupied table}) \propto \text{number of customers at table}$
  - $P(\text{next unoccupied table}) \propto \alpha$

## Chinese Restaurant Process (4)



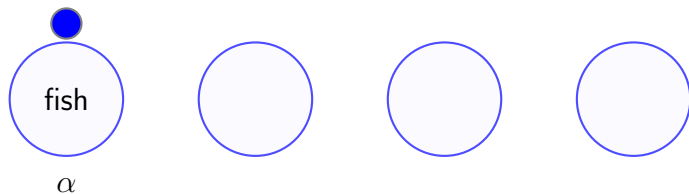
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2, 1$
- $P(\mathbf{z}) = \alpha/\alpha \times 1/(1 + \alpha) \times \alpha/(2 + \alpha) \times 2/(3 + \alpha)$
- Next customer chooses a table according to:
  - $P(\text{occupied table}) \propto \text{number of customers at table}$
  - $P(\text{next unoccupied table}) \propto \alpha$

# Labelled Chinese Restaurant Process (0)



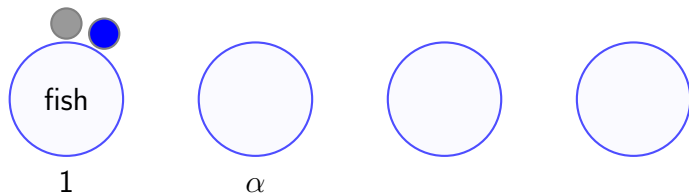
- Table  $\rightarrow$  label mapping  $\mathbf{y} =$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} =$
- Output sequence  $\mathbf{x} =$
- $P(\mathbf{x}) = 1$
  
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i = k$ ) share label  $y_k$
- Customer  $i$  sitting at table  $z_i$  has label  $x_i = y_{z_i}$

# Labelled Chinese Restaurant Process (1)



- Table  $\rightarrow$  label mapping  $\mathbf{y} = \text{fish}$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1$
- Output sequence  $\mathbf{x} = \text{fish}$
- $P(\mathbf{x}) = \alpha/\alpha \times P_0(\text{fish})$
  
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i = k$ ) share label  $y_k$
- Customer  $i$  sitting at table  $z_i$  has label  $x_i = y_{z_i}$

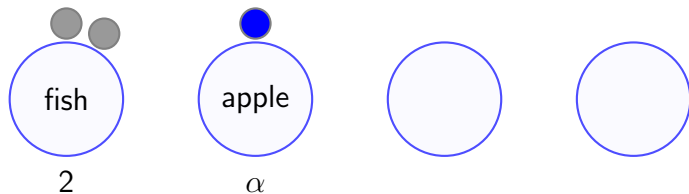
## Labelled Chinese Restaurant Process (2)



- Table  $\rightarrow$  label mapping  $\mathbf{y} = \text{fish}$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1$
- Output sequence  $\mathbf{x} = \text{fish}, \text{fish}$
- $P(\mathbf{x}) = P_0(\text{fish}) \times 1/(1 + \alpha)$
  
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i = k$ ) share label  $y_k$
- Customer  $i$  sitting at table  $z_i$  has label  $x_i = y_{z_i}$

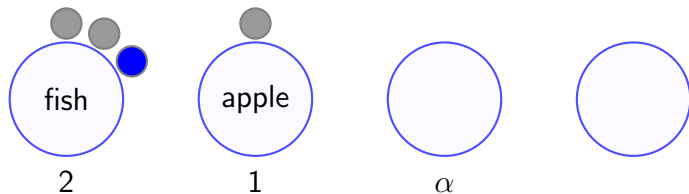


## Labelled Chinese Restaurant Process (3)



- Table  $\rightarrow$  label mapping  $\mathbf{y} = \text{fish, apple}$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2$
- Output sequence  $\mathbf{x} = \text{fish, fish, apple}$
- $P(\mathbf{x}) = P_0(\text{fish}) \times 1/(1 + \alpha) \times \alpha/(2 + \alpha)P_0(\text{apple})$
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i = k$ ) share label  $y_k$
- Customer  $i$  sitting at table  $z_i$  has label  $x_i = y_{z_i}$

## Labelled Chinese Restaurant Process (4)



- Table  $\rightarrow$  label mapping  $\mathbf{y} = \text{fish, apple}$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2$
- Output sequence  $\mathbf{x} = \text{fish, fish, apple, fish}$
- $P(\mathbf{x}) = P_0(\text{fish}) \times 1/(1 + \alpha) \times \alpha/(2 + \alpha) P_0(\text{apple}) \times 2/(3 + \alpha)$
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i = k$ ) share label  $y_k$
- Customer  $i$  sitting at table  $z_i$  has label  $x_i = y_{z_i}$

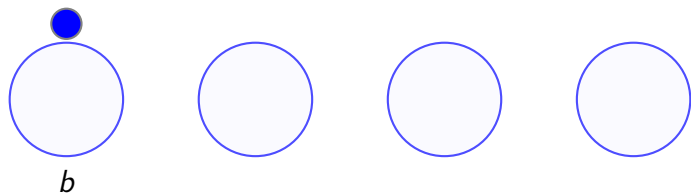
# Pitman-Yor Process (0)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} =$
- $P(\mathbf{z}) = 1$
  
- In CRPs, probability of choosing a table  $\propto$  number of customers  
 $\Rightarrow$  strong *rich get richer* effect
- Pitman-Yor processes take mass  $a$  from each occupied table and give it to the new table

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) - a & \text{if } k \leq m = \max(\mathbf{z}) \\ ma + b & \text{if } k = m + 1 \end{cases}$$

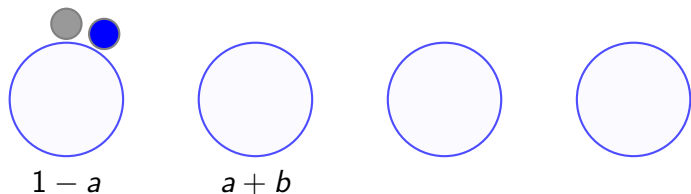
# Pitman-Yor Process (1)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1$
- $P(\mathbf{z}) = b/b$
- In CRPs, probability of choosing a table  $\propto$  number of customers  
 $\Rightarrow$  strong *rich get richer* effect
- Pitman-Yor processes take mass  $a$  from each occupied table and give it to the new table

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) - a & \text{if } k \leq m = \max(\mathbf{z}) \\ ma + b & \text{if } k = m + 1 \end{cases}$$

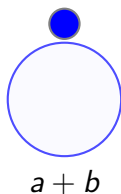
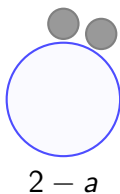
## Pitman-Yor Process (2)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1$
- $P(\mathbf{z}) = b/b \times (1 - a)/(1 + b)$
- In CRPs, probability of choosing a table  $\propto$  number of customers  
 $\Rightarrow$  strong *rich get richer* effect
- Pitman-Yor processes take mass  $a$  from each occupied table and give it to the new table

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) - a & \text{if } k \leq m = \max(\mathbf{z}) \\ ma + b & \text{if } k = m + 1 \end{cases}$$

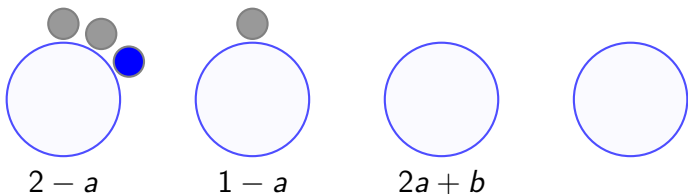
## Pitman-Yor Process (3)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2$
- $P(\mathbf{z}) = b/b \times (1 - a)/(1 + b) \times (a + b)/(2 + b)$
- In CRPs, probability of choosing a table  $\propto$  number of customers  
 $\Rightarrow$  strong *rich get richer* effect
- Pitman-Yor processes take mass  $a$  from each occupied table and give it to the new table

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) - a & \text{if } k \leq m = \max(\mathbf{z}) \\ ma + b & \text{if } k = m + 1 \end{cases}$$

## Pitman-Yor Process (4)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2, 1$
- $P(\mathbf{z}) = b/b \times (1 - a)/(1 + b) \times (a + b)/(2 + b) \times (2 - a)/(3 + b)$
- In CRPs, probability of choosing a table  $\propto$  number of customers  
 $\Rightarrow$  strong *rich get richer* effect
- Pitman-Yor processes take mass  $a$  from each occupied table and give it to the new table

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) - a & \text{if } k \leq m = \max(\mathbf{z}) \\ ma + b & \text{if } k = m + 1 \end{cases}$$

# Summary: Chinese Restaurant Processes and Pitman-Yor Processes

- *Chinese Restaurant Processes* (CRPs) generalize Dirichlet-Multinomials to an *unbounded number of outcomes*
  - ▶ *concentration parameter*  $\alpha$  controls how likely a new outcome is
  - ▶ CRPs exhibit a *rich get richer* power-law behaviour
- *Pitman-Yor Processes* (PYPs) generalize CRPs by adding an additional parameter (each PYP has  $a$  and  $b$  parameters)
  - ▶ PYPs can describe a wider range of distributions than CRPs
- *Labelled CRPs and PYPs* use a *base distribution* to label each table
  - ▶ base distribution can have *infinite support*
  - ▶ concentrates mass on a countable subset



# Outline

Non-parametric inference in word segmentation

Chinese Restaurant Processes

Adaptor grammars

- Adaptor grammars for word segmentation

- Priors on Adaptor Grammar PYP parameters

Unsupervised inference of Adaptor Grammars

- Modal word segmentation

- Random vs incremental initialization

- Table label resampling

Conclusion

# Adaptor grammars: informal description

- The trees generated by an adaptor grammar are defined by CFG rules as in a CFG
- A subset of the nonterminals are *adapted*
- *Unadapted nonterminals* expand by picking a rule and recursively expanding its children, as in a PCFG
- *Adapted nonterminals* can expand in two ways:
  - ▶ by picking a rule and recursively expanding its children, or
  - ▶ by generating a previously generated tree (with probability proportional to the number of times previously generated)
- Implemented by having a CRP/PYP for each adapted nonterminal
- The CFG rules of the adapted nonterminals determine the *base distributions* of these CRPs/PYPs

# Adaptor grammar for stem-suffix morphology (0)

Word  $\rightarrow$  Stem Suffix



Stem  $\rightarrow$  Phoneme<sup>+</sup>



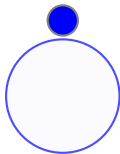
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



Generated words:

# Adaptor grammar for stem-suffix morphology (1a)

Word  $\rightarrow$  Stem Suffix



Stem  $\rightarrow$  Phoneme<sup>+</sup>



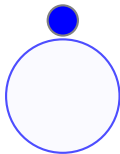
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



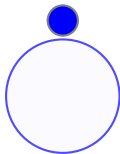
Generated words:

# Adaptor grammar for stem-suffix morphology (1b)

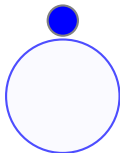
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



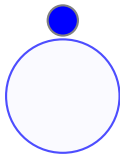
Suffix → Phoneme<sup>\*</sup>



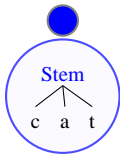
Generated words:

# Adaptor grammar for stem-suffix morphology (1c)

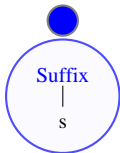
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



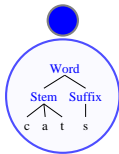
Suffix → Phoneme<sup>\*</sup>



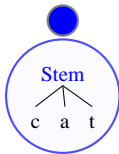
Generated words:

# Adaptor grammar for stem-suffix morphology (1d)

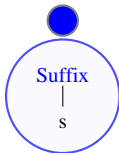
Word  $\rightarrow$  Stem Suffix



Stem  $\rightarrow$  Phoneme<sup>+</sup>



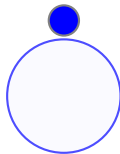
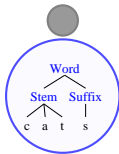
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



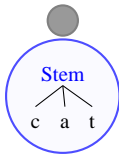
Generated words: **cats**

# Adaptor grammar for stem-suffix morphology (2a)

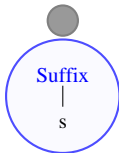
Word  $\rightarrow$  Stem Suffix



Stem  $\rightarrow$  Phoneme<sup>+</sup>



Suffix  $\rightarrow$  Phoneme<sup>\*</sup>

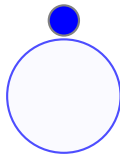
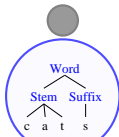


Generated words: cats

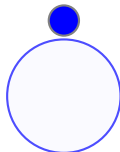
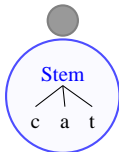


# Adaptor grammar for stem-suffix morphology (2b)

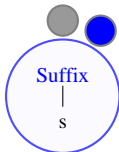
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



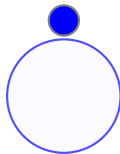
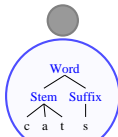
Suffix → Phoneme<sup>\*</sup>



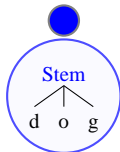
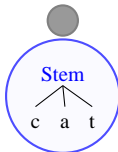
Generated words: cats

# Adaptor grammar for stem-suffix morphology (2c)

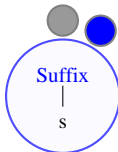
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



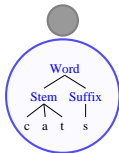
Suffix → Phoneme<sup>\*</sup>



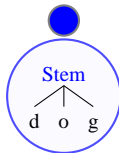
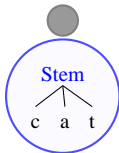
Generated words: cats

# Adaptor grammar for stem-suffix morphology (2d)

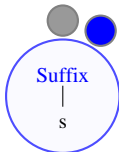
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



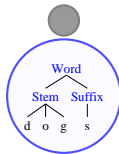
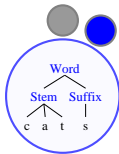
Suffix → Phoneme<sup>\*</sup>



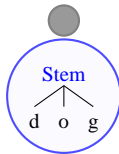
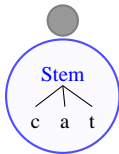
Generated words: cats, dogs

# Adaptor grammar for stem-suffix morphology (3)

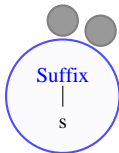
Word  $\rightarrow$  Stem Suffix



Stem  $\rightarrow$  Phoneme<sup>+</sup>



Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



Generated words: cats, dogs, **cats**

# Adaptor grammars as generative processes

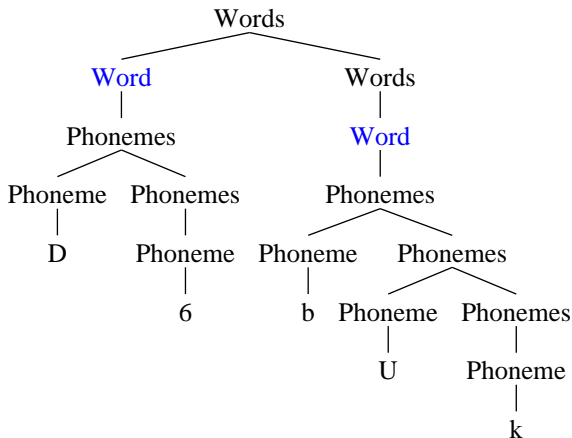
- The sequence of trees generated by an adaptor grammar are *not* independent
  - ▶ it *learns* from the trees it generates
  - ▶ if an adapted subtree has been used frequently in the past, it's more likely to be used again
- but the sequence of trees is *exchangable* (important for sampling)
- An *unadapted nonterminal*  $A$  expands using  $A \rightarrow \beta$  with probability  $\theta_{A \rightarrow \beta}$
- Each adapted nonterminal  $A$  is associated with a CRP (or PYP) that caches previously generated subtrees rooted in  $A$
- An *adapted nonterminal*  $A$  expands:
  - ▶ to a subtree  $\tau$  rooted in  $A$  with probability proportional to the number of times  $\tau$  was previously generated
  - ▶ using  $A \rightarrow \beta$  with probability proportional to  $\alpha_A \theta_{A \rightarrow \beta}$

# Unigram word segmentation adaptor grammar

Sentence  $\rightarrow$  Word<sup>+</sup>

Word  $\rightarrow$  Phoneme<sup>+</sup>

- *Adapted nonterminals* indicated by underlining



- Adapting Words means that the grammar learns the probability of each Word subtree independently
- Unigram word segmentation on Brent corpus: 56% token f-score

# Unigram adaptor grammar after learning

- Given the Brent corpus and the unigram adaptor grammar

Words  $\rightarrow$  Word<sup>+</sup>

Word  $\rightarrow$  Phon<sup>+</sup>

the learnt adapted grammar contains 1,712 rules such as:

15758 Words  $\rightarrow$ Word Words

9791 Words  $\rightarrow$ Word

1660 Word  $\rightarrow$  Phon<sup>+</sup>

402 Word  $\rightarrow$  *y u*

137 Word  $\rightarrow$  *l n*

111 Word  $\rightarrow$  *w l T*

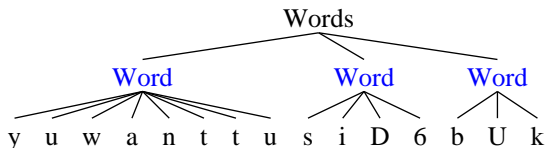
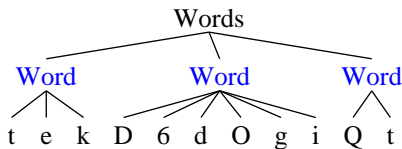
100 Word  $\rightarrow$  *D 6 d O g i*

45 Word  $\rightarrow$  *l n D 6*

20 Word  $\rightarrow$  *l n D 6 h Q s*

## unigram: Words

- Unigram word segmentation model assumes each word is generated independently
- But there are strong inter-word dependencies (collocations)
- Unigram model can only capture such dependencies by analyzing collocations as words (Goldwater 2006)



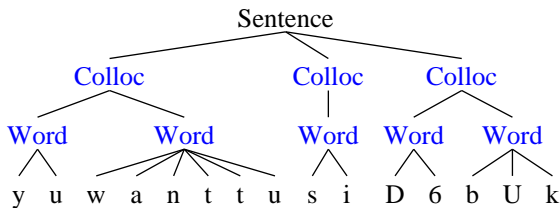


# colloc: Collocations $\Rightarrow$ Words

Sentence  $\rightarrow$  Colloc<sup>+</sup>

Colloc  $\rightarrow$  Word<sup>+</sup>

Word  $\rightarrow$  Phon<sup>+</sup>



- A Colloc(ation) consists of one or more words
- Both Words and Collocs are adapted (learnt)
- Significantly improves word segmentation accuracy over unigram model (76% f-score;  $\approx$  Goldwater's bigram model)

# colloc-syll: Collocations $\Rightarrow$ Words $\Rightarrow$ Syllables

Sentence  $\rightarrow$  Colloc<sup>+</sup>

Word  $\rightarrow$  SyllableIF

Word  $\rightarrow$  SyllableI (Syllable) (Syllable) SyllableF

Onset  $\rightarrow$  Consonant<sup>+</sup>

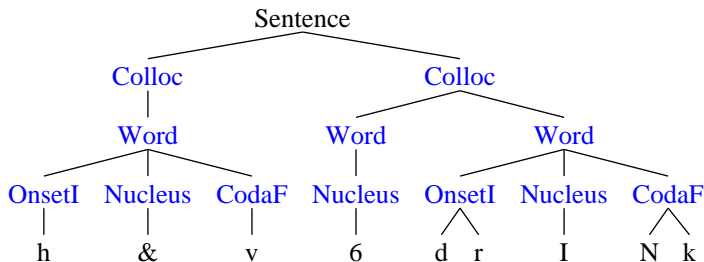
Nucleus  $\rightarrow$  Vowel<sup>+</sup>

Colloc  $\rightarrow$  Word<sup>+</sup>

Syllable  $\rightarrow$  (Onset) Rhyme

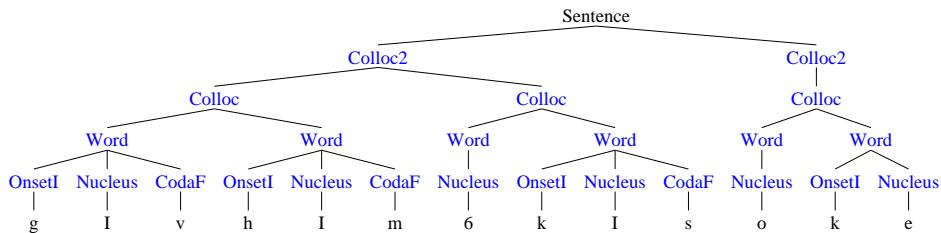
Rhyme  $\rightarrow$  Nucleus (Coda)

Coda  $\rightarrow$  Consonant<sup>+</sup>



- With 2 Collocation levels, f-score = 87%

# colloc-syll: Collocations $\Rightarrow$ Words $\Rightarrow$ Syllables



# Another application of adaptor grammars: Learning structure in names

- Many different kinds of names
  - ▶ Person names, e.g., *Mr. Sam Spade Jr.*
  - ▶ Company names, e.g., *United Motor Manufacturing Corp.*
  - ▶ Other names, e.g., *United States of America*
- At least some of these are structured; e.g., *Mr* is an honorific, *Sam* is first name, *Spade* is a surname, etc.
- Penn treebanks assign flat structures to base NPs (including names)
- Data set: 10,787 unique lowercased sequences of base NP proper nouns, containing 23,392 words
- Can we automatically learn the structure of these names?

# Adaptor grammar for names

NP  $\rightarrow$  Unordered<sup>+</sup>

NP  $\rightarrow$  (A0) (A1) ... (A6)

A0  $\rightarrow$  Word<sup>+</sup>

...

A6  $\rightarrow$  Word<sup>+</sup>

Unordered  $\rightarrow$  Word<sup>+</sup>

NP  $\rightarrow$  (B0) (B1) ... (B6)

B0  $\rightarrow$  Word<sup>+</sup>

...

B6  $\rightarrow$  Word<sup>+</sup>

- *Warning: hand-selected output, no evaluation!*

(A0 barrett) (A3 smith)

(A0 albert) (A2 j.) (A3 smith) (A4 jr.)

(A0 robert) (A2 b.) (A3 van dover)

(B0 aim) (B1 prime rate) (B2 plus) (B5 fund) (B6 inc.)

(B0 balfour) (B1 maclaine) (B5 international) (B6 ltd.)

(B0 american express) (B1 information services) (B6 co)

(U abc) (U sports)

(U sports illustrated)

(U sports unlimited)

# Bayesian inference for PYP parameters

- Adaptor grammars have 1 (CRP) or 2 (PYP) hyper-parameters for each adapted non-terminal  $X$
- Previous work used CRP adaptors with *tied parameters*
- Bayesian prior: for each adapted nonterminal  $X$

$$a_X \sim \text{Beta}(1, 1)$$

$$b_X \sim \text{Gamma}(10, 0.1)$$

- ▶  $\text{Gamma}(10, 0.1)$  is a vague Gamma prior (MacKay 2003)
- ▶ permits  $a_X$  and  $b_X$  to *vary with adapted nonterminal  $X$*
- Estimate with *slice sampling* (no proposal distribution; Neal 2003)
- *Biggest improvement on complex models*, e.g., colloc-syll:
  - ▶ tied parameters 78% token f-score
  - ▶  $a_X = 0$ , sampling  $b_X$  (i.e., CRP adaptors) 84% token f-score
  - ▶ sampling  $a_X$  and  $b_X$  (i.e., PYP adaptors) 87% token f-score

# Outline

Non-parametric inference in word segmentation

Chinese Restaurant Processes

Adaptor grammars

- Adaptor grammars for word segmentation

- Priors on Adaptor Grammar PYP parameters

Unsupervised inference of Adaptor Grammars

- Modal word segmentation

- Random vs incremental initialization

- Table label resampling

Conclusion

# Unsupervised inference via Gibbs sampling

- Observations (terminal strings)  $\mathbf{x} = (x_1, \dots, x_n)$   
Hidden labels (parse trees)  $\mathbf{t} = (t_1, \dots, t_n)$   
Probabilistic model (adaptor grammar)  $P(\mathbf{x}, \mathbf{t})$
- Gibbs sampling algorithm:
  - initialize  $\mathbf{t}$  somehow (e.g., random trees)
  - repeat forever:
    - pick an index  $j \in 1, \dots, n$  at random
    - replace  $t_j$  with a random sample from  $P(t \mid x_j, \mathbf{t}_{-j})$   
where  $\mathbf{t}_{-j} = (t_1, \dots, t_{j-1}, t_{j+1}, \dots, t_n)$
- After *burn-in* the samples  $\mathbf{t}$  are distributed according to  $P(\mathbf{t} \mid \mathbf{x})$



# Finding the modal word segmentation

- Previous work decoded using last sampled trees (2,000 epochs)
- After burn-in, samples are distributed according to  $P(\mathbf{t} \mid \mathbf{x})$   
⇒ use samples to identify *modal word segmentation*
- Modal decoding:
  - ▶ For each sentence  $x_i$  collect *sample parses*  $\mathbf{s}_i = (s_i^{(1)}, \dots, s_i^{(800)})$   
(every 10th epoch from epochs 1,000–2,000 from 8 runs)
  - ▶ Compute *word segmentations*  $\mathbf{w}_i = (w_i^{(1)}, \dots, w_i^{(800)})$  from parses
  - ▶ Compute *modal segmentation*  $\hat{w}_i = \operatorname{argmax}_w n_w(\mathbf{w}_i)$ ,  
where  $n_w(\mathbf{w}_i)$  is the number of times  $w$  appears in  $\mathbf{w}_i$
- Improves word segmentation token f-score in all models

Model	Average	Max-modal
unigram	55%	56%
colloc	74%	76%
colloc-syll	85%	87%

- Goodman (1998) max-marginal decoding should also be possible

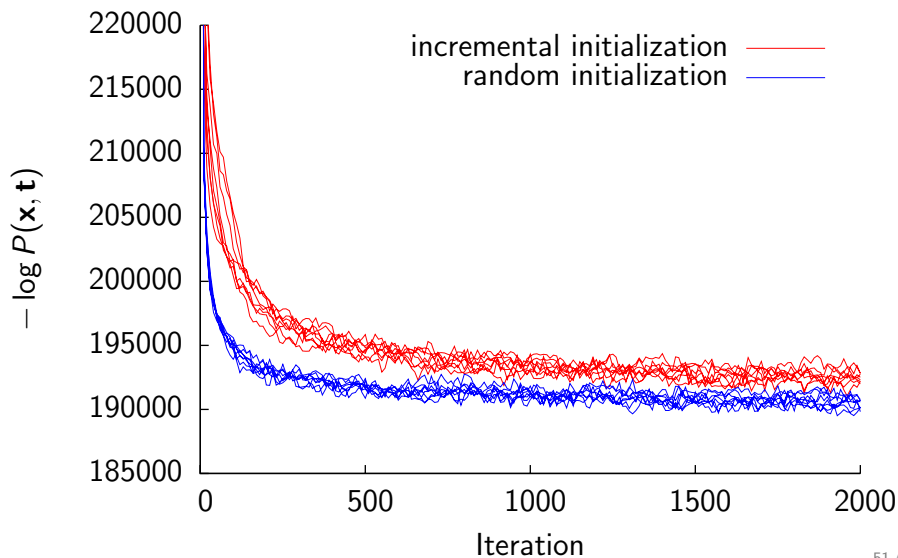
# Random vs incremental initialization

- The Gibbs sampler parse trees  $\mathbf{t}$  needs to be initialized somehow  
**Random initialization:** Assign each string  $x_i$  a random parse  $t_i$  generated by base PCFG  
**Incremental initialization:** Sample  $t_i$  from  $P(t \mid x_i, \mathbf{t}_{1:i-1})$
- Incremental initialization is easy to implement in a Gibbs sampler
- Incremental initialization improves token f-score in all models, especially on simple models

Model	Random	Incremental
unigram	56%	81%
colloc	76%	86%
colloc-syll	87%	89%

*but see caveats on next slide!*

# Incremental initialization produces low-probability parses



# Why incremental initialization produces low-probability parses

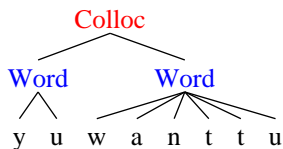
- Incremental initialization produces sample parses  $\mathbf{t}$  with lower probability  $P(\mathbf{t} \mid \mathbf{x})$
- Possible explanation: (Goldwater's 2006 analysis of Brent's model)
  - ▶ All the models tend to *undersegment* (i.e., find collocations instead of words)
  - ▶ Incremental initialization *greedily searches for common substrings*
  - ▶ Shorter strings are more likely to be recur early than longer ones

# Table label resampling

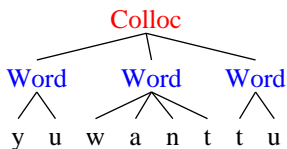
- Each adapted non-terminal has a CRP with tables labelled with parses
- “Rich get richer”  $\Rightarrow$  resampling a sentence’s parse reuses the same cached subtrees
- *Resample table labels* as well sentence parses
  - ▶ A table label may be used in many sentence parses
  - $\Rightarrow$  Resampling a single table label may change the parses of a single sentence
  - $\Rightarrow$  table label resampling can improve mobility with grammars with a hierarchy of adapted non-terminals
- Essential for grammars with a complex hierarchical structure

# Table label resampling example

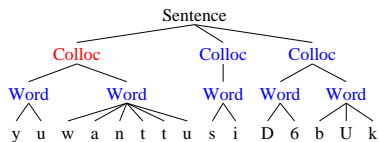
Label on table in Chinese Restaurant for colloc



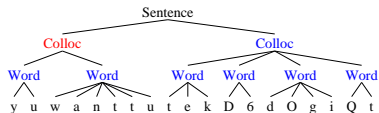
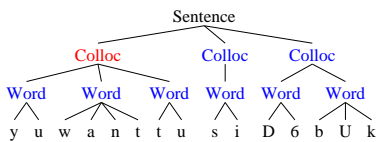
⇒



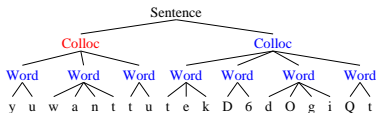
Resulting changes in parse trees



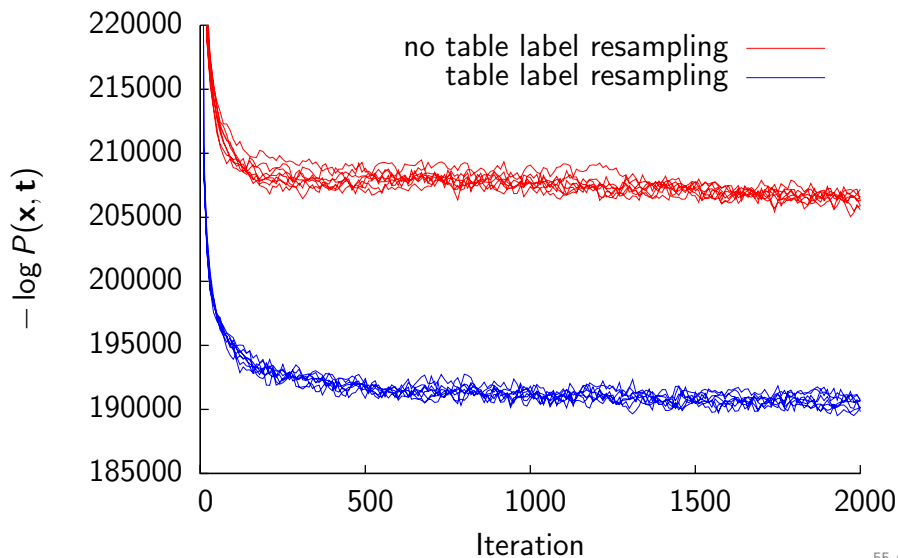
⇒



⇒



## Table label resampling produces much higher-probability parses



# Outline

Non-parametric inference in word segmentation

Chinese Restaurant Processes

Adaptor grammars

- Adaptor grammars for word segmentation

- Priors on Adaptor Grammar PYP parameters

Unsupervised inference of Adaptor Grammars

- Modal word segmentation

- Random vs incremental initialization

- Table label resampling

Conclusion



# Conclusion

- MCMC sampling is a natural approach to inference in non-parametric models
- Put Bayesian priors on PYP hyperparameters and slice sample
  - ▶ improves colloc-syll model by 9%
- Modal decoding from multiple samples
  - ▶ improves colloc-syll model by 2%
- Random initialization instead of incremental initialization
  - ▶ hurts colloc-syll model by 2%, but produces higher probability parses
- Table label resampling in hierarchical CRPs/PYPs
  - ▶ improves colloc-syll model by 20%
- Taken together, performance of colloc-syll model improves from last year's 78% to 87% token f-score.