

# Learning rules with Adaptor Grammars

Mark Johnson

joint work with Sharon Goldwater and Tom Griffiths

May, 2008

# The drunk under the lamppost

Late one night, a drunk guy is crawling around under a lamppost. A cop comes up and asks him what he's doing.

*“I'm looking for my keys,”* the drunk says. *“I lost them about three blocks away.”*

*“So why aren't you looking for them where you dropped them?”* the cop asks.

The drunk looks at the cop, amazed that he'd ask so obvious a question. *“Because the light is better here.”*

# Ideas behind talk

- Most successful statistical learning methods are *parametric*
  - ▶ PCFGs have one probability parameter per rule
  - ▶ PCFG learning: given rules and data, learn rule probabilities
- Non-parametric learning: learn parameters (rules) as well as values
- *Adaptor grammars*:
  - ▶ are a framework for specifying hierarchical nonparametric Bayesian models
  - ▶ can express a variety of linguistically-interesting structures
  - ▶ are approximated by PCFGs, where number of rules depends on data

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

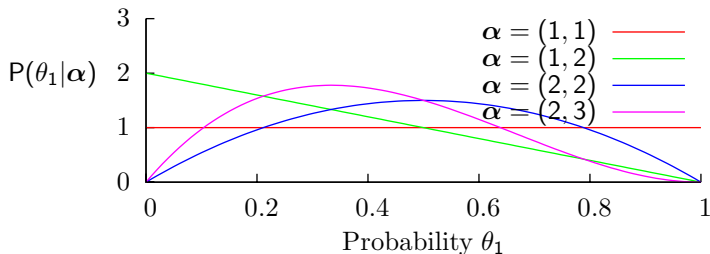
# Language acquisition as Bayesian inference

$$\underbrace{P(\text{Grammar} \mid \text{Data})}_{\text{Posterior}} \propto \underbrace{P(\text{Data} \mid \text{Grammar})}_{\text{Likelihood}} \underbrace{P(\text{Grammar})}_{\text{Prior}}$$

- Likelihood measures how well grammar describes data
- Prior expresses knowledge of grammar before data is seen
  - ▶ can be very specific (e.g., Universal Grammar)
  - ▶ can be very general (e.g., prefer shorter grammars)
- Posterior is *distribution* over grammars
  - ▶ expresses uncertainty about which grammar is correct
- But: *infinitely many* grammars may be consistent with Data

# Multinomial distributions and Dirichlet priors

- $X$  follows a *multinomial distribution* if  $X$  ranges over  $1, \dots, m$  and  $P(X = k|\boldsymbol{\theta}) = \theta_k$
- Bayesian inference for  $\boldsymbol{\theta}$ :  $P(\boldsymbol{\theta}|X) \propto P(X|\boldsymbol{\theta})P(\boldsymbol{\theta})$
- A *Dirichlet distribution*  $P(\boldsymbol{\theta}|\boldsymbol{\alpha})$  is a probability distribution over multinomial parameters  $\boldsymbol{\theta}$ 
  - ▶ One Dirichlet parameter  $\alpha_k$  for each outcome  $k$
- Dirichlets are *conjugate prior* for multinomials
  - ▶ If prior is Dirichlet  $(\alpha_1, \dots, \alpha_k, \dots, \alpha_m)$  and we observe data  $X = k$ , posterior is Dirichlet  $(\alpha_1, \dots, \alpha_k + 1, \dots, \alpha_m)$



## Conditional Dirichlet-multinomial distribution

- Suppose we don't know  $\boldsymbol{\theta}$ , but only know Dirichlet prior  $\boldsymbol{\alpha}$ :

$$P(X = k \mid \boldsymbol{\alpha}) = \int P(X = k \mid \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) d\boldsymbol{\theta} \propto \alpha_k$$

- Suppose we observe  $\mathbf{X} = (x_1, \dots, x_n)$ . What is probability of next outcome  $X_{n+1}$ ?

$$P(X_{n+1} = k \mid \boldsymbol{\alpha}, \mathbf{X}) \propto \alpha_k + n_k(\mathbf{X})$$

where  $n_k(\mathbf{X})$  is number of times  $k$  appears in  $\mathbf{X}$

- Example: coin with sides  $h, t$ . Prior  $\alpha_h = \alpha_t = 1$ .  $\mathbf{X} = (t, h, t)$ . Posterior  $\alpha'_h = 2, \alpha'_t = 3$ , so  $P(X_4 = h \mid \boldsymbol{\alpha}, \mathbf{X} = (t, h, t)) = 2/5$

# Predicting a sequence of outcomes

- Given Dirichlet prior  $\boldsymbol{\alpha}$  and observations  $\mathbf{X}$ , what is probability of next *two* outcomes  $X_{n+1}$  and  $X_{n+2}$ ?

$$\begin{aligned} & \mathbb{P}(X_{n+1} = k, X_{n+2} = k' \mid \boldsymbol{\alpha}, \mathbf{X}) \\ &= \mathbb{P}(X_{n+1} = k \mid \boldsymbol{\alpha}, \mathbf{X}) \mathbb{P}(X_{n+2} = k' \mid \boldsymbol{\alpha}, \mathbf{X}, X_{n+1} = k) \\ &\propto (\alpha_k + n_k(\mathbf{X})) (\alpha_{k'} + n_{k'}(\mathbf{X}) + \delta_{k,k'}) \end{aligned}$$

i.e., the probability that  $X_{n+2} = k'$  is affected by value of  $X_{n+1}$

$\Rightarrow$  *The probability of each outcome type changes as we progress through a sequence of outcomes*



# Outline

Bayesian inference of multinomials

**Inference of Probabilistic Context-Free Grammars**

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

# Probabilistic context-free grammars

- Rules in *Context-Free Grammars* (CFGs) expand nonterminals into sequences of terminals and nonterminals
- A *Probabilistic CFG* (PCFG) associates each nonterminal with a multinomial distribution over the rules that expand it
- Probability of a tree is the *product of the probabilities of the rules* used to construct it

Rule $r$	$\theta_r$
$S \rightarrow NP VP$	1.0
$NP \rightarrow \text{Sam}$	0.75
$VP \rightarrow \text{barks}$	0.6

Rule $r$	$\theta_r$
$NP \rightarrow \text{Sandy}$	0.25
$VP \rightarrow \text{snores}$	0.4

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{Sam} \quad \text{barks} \end{array} \right) = 0.45$$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{Sandy} \quad \text{snores} \end{array} \right) = 0.1$$

## Dirichlet priors on rule probabilities

- Prior  $P(\boldsymbol{\theta})$  specifies probability of rule probabilities  $\boldsymbol{\theta}$
- Conjugate prior is a *product of Dirichlets* (one for each nonterminal)
- Prior has a parameter  $\alpha_{A \rightarrow \beta}$  for each rule  $A \rightarrow \beta$
- Suppose we're given prior  $\boldsymbol{\alpha}$  and observe parse trees  $\mathbf{T} = (T_1, \dots, T_n)$ . What is posterior  $P(\boldsymbol{\theta} \mid \boldsymbol{\alpha}, \mathbf{T})$ ?
- $P(\boldsymbol{\theta} \mid \boldsymbol{\alpha}, \mathbf{T})$  is a product of Dirichlets with parameters  $\boldsymbol{\alpha}'$ , where:

$$\alpha'_{A \rightarrow \beta} = \alpha_{A \rightarrow \beta} + n_{A \rightarrow \beta}(\mathbf{T})$$

and  $n_{A \rightarrow \beta}(\mathbf{T})$  is number of times  $A \rightarrow \beta$  in  $\mathbf{T}$

In summary:  $\boldsymbol{\alpha}' = \boldsymbol{\alpha} + \mathbf{n}(\mathbf{T})$

⇒ Bayesian estimation is easy given visible data (parse trees) and Dirichlet prior

# Estimating $\theta$ from strings via Gibbs sampling

- Learning from terminal strings  $\mathbf{W} \Rightarrow$  parse trees  $\mathbf{T}$  are *hidden*
- No known closed form for posterior  $P(\theta|\alpha, \mathbf{W})$ , but we can approximate it by sampling
- *Gibbs sampling* is a Markov Chain Monte Carlo (MCMC) method for sampling from  $P(X_1, \dots, X_n)$

initialize  $x_1, \dots, x_n$  somehow

repeat forever:

for  $i$  in  $1, \dots, n$ :

set  $x_i$  to a sample from  $P(X_i|\mathbf{X}_{-i} = \mathbf{x}_{-i})$

where  $\mathbf{X}_{-i} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ .

Under very general conditions, after “burn-in” this sequence of samples has distribution  $P(X_1, \dots, X_n)$

- Example: to produce samples from  $P(X, Y)$ , guess  $x, y$ , sample  $X^{(1)}$  from  $P(X|Y^{(0)})$ , then sample  $Y^{(1)}$  from  $P(Y|X)$ , then  $X$  from  $P(X|Y)$ , then  $Y$  from  $P(Y|X)$ , ...

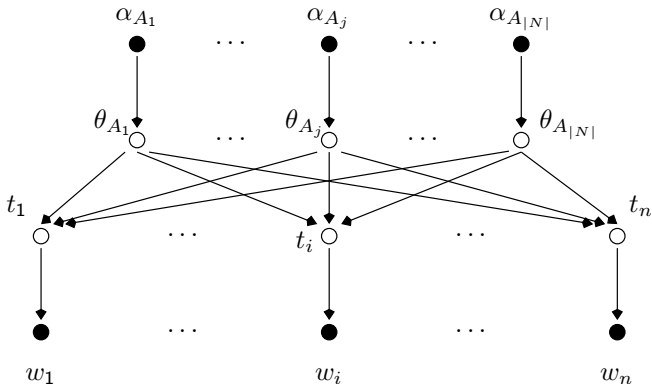
# Bayes net representation of Bayesian PCFG

*Prior  $\alpha$*

*Rule probabilities  $\theta$*

*Parse trees  $T$*

*Terminal strings  $W$*



$$\begin{array}{l|l}
 \theta_A & \alpha_A \sim \text{Dir}(\alpha_A) \\
 T_i & \theta \sim G_S(\theta) \\
 W_i & T_i = \text{YIELD}(T_i)
 \end{array}
 \quad
 \begin{array}{l}
 (\theta_A \text{ are prob. of rules expanding } A) \\
 (T_i \text{ is parse tree for string } W_i)
 \end{array}$$

# A Gibbs sampler for trees and rule probabilities

- Generative model:

$$P(\mathbf{W}, \mathbf{T}, \boldsymbol{\theta} \mid \boldsymbol{\alpha}) = P(\mathbf{W} \mid \mathbf{T}) P(\mathbf{T} \mid \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \boldsymbol{\alpha})$$

- Dirichlet prior parameters  $\boldsymbol{\alpha}$  and strings  $\mathbf{W}$  are given
- Rule probabilities  $\boldsymbol{\theta}$  and parse trees  $\mathbf{T}$  are unknown
- Goal: estimate  $P(\mathbf{T}, \boldsymbol{\theta} \mid \boldsymbol{\alpha}, \mathbf{W})$
- Gibbs sampler iteration:
  - ▶ sample  $\boldsymbol{\theta}$  from  $P(\boldsymbol{\theta} \mid \mathbf{T}, \boldsymbol{\alpha})$ 
    - i.e., sampling multinomials from Dirichlets
  - ▶ sample  $\mathbf{T}$  from  $P(\mathbf{T} \mid \mathbf{W}, \boldsymbol{\theta})$ 
    - sampling parse tree  $T_i$  for string  $W_i$  given  $\boldsymbol{\theta}$
    - can be done in parallel for each string  $W_i$
    - $O(|W_i|^3)$  dynamic programming algorithm

## Integrating out the rule probabilities

- Posterior distribution  $P(\boldsymbol{\theta} \mid \mathbf{T}, \boldsymbol{\alpha})$  is Dirichlet with parameters  $\boldsymbol{\alpha}' = \boldsymbol{\alpha} + \mathbf{n}(\mathbf{T})$

⇒ Can calculate conditional probability of next parse tree  $T_n$  given previous parses  $\mathbf{T}_{-n} = (T_1, \dots, T_{n-1})$  and  $\boldsymbol{\alpha}$

$$\begin{aligned} P(T_n \mid \mathbf{T}_{-n}, \boldsymbol{\alpha}) &= \int P(T_n \mid \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathbf{T}_{-n}, \boldsymbol{\alpha}) d\boldsymbol{\theta} \\ &= P(T_n \mid \boldsymbol{\alpha}') \quad \text{where } \boldsymbol{\alpha}' = \boldsymbol{\alpha} + \mathbf{n}(\mathbf{T}_{-n}) \end{aligned}$$

⇒ Gibbs sampler with parse trees  $T_i$  as the  $n$  components

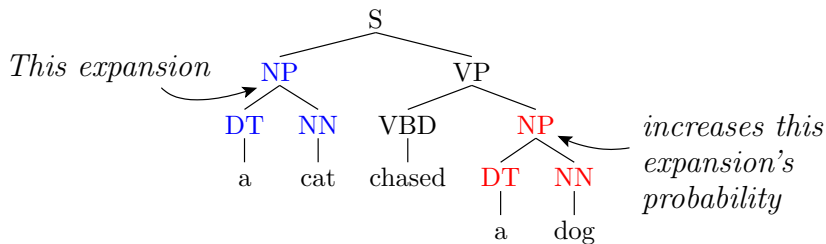
- ▶ requires sampling  $T_i$  from  $P(T_i \mid W_i, \mathbf{T}_{-i}, \boldsymbol{\alpha})$
- ▶ but calculating normalizing constant in

$$P(T_i \mid W_i, \mathbf{T}_{-i}, \boldsymbol{\alpha}) = \frac{P(T_i \mid \mathbf{T}_{-i}, \boldsymbol{\alpha})}{\sum_{T: \text{YIELD}(T)=W_i} P(T \mid \mathbf{T}_{-i}, \boldsymbol{\alpha})}$$

is hard!

# Conditional parse probability in collapsed model

- Dirichlet prior on  $\theta$ 
  - $\Rightarrow$  probability of each rule expansion  $A \rightarrow \beta \propto \alpha_{A \rightarrow \beta}$
- Parse tree is a *sequence* of rule expansions
  - $\Rightarrow$  Update Dirichlet prior after each rule expansion
  - $\Rightarrow$   $A$ 's expansion probabilities *change* each time  $A$  expands
  - $\Rightarrow$  A rule's expansion "primes" its subsequent use



- A rule's probability depends on which rules were used earlier
  - $\Rightarrow$  *no dynamic programming algorithms*
  - $\Rightarrow$  no efficient way of calculating normalizing constant



# A Metropolis-within-Gibbs sampler for trees

- Metropolis “accept/reject” sampling procedures only require “unnormalized” probabilities
- But they require a *proposal distribution* that is
  - ▶ easy to sample from, and
  - ▶ not too different from true distribution
- Proposal distribution for  $P(T_i \mid \alpha, \mathbf{T}_{-i}, W_i)$  is PCFG with rule probabilities

$$\theta'_{A \rightarrow \beta} \propto \alpha_{A \rightarrow \beta} + n_{A \rightarrow \beta}(\mathbf{T}_{-i})$$

- Metropolis-within-Gibbs sampling algorithm:

repeat forever:

for  $i$  in  $1, \dots, n$ :

compute proposal PCFG prob.  $\theta'$  from  $\alpha + n(\mathbf{T}_{-i})$

sample parse tree  $T^*$  from  $P(T \mid W_i, \theta')$

compute “true” probability  $P(T^* \mid \alpha, \mathbf{T}_{-i})$

Metropolis accept/reject  $T^*$  as new value for  $T_i$

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

**Limitations of PCFG learning**

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

# Learning syntactic structure is hard

- Bayesian PCFG estimation works well on toy data
- Results are disappointing on “real” data
  - ▶ wrong rules?
  - ▶ rules need to be given a priori; can we learn them too?
- Strategy: study simpler cases
  - ▶ Morphological segmentation (e.g., *walking* = *walk+ing*)
  - ▶ Word segmentation of unsegmented utterances

# A CFG for stem-suffix morphology

Word  $\rightarrow$  Stem Suffix

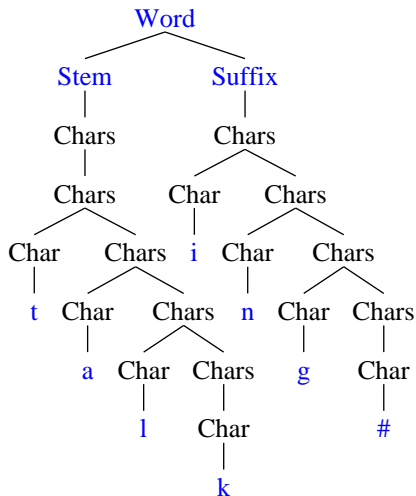
Stem  $\rightarrow$  Chars

Suffix  $\rightarrow$  Chars

Chars  $\rightarrow$  Char

Chars  $\rightarrow$  Char Chars

Char  $\rightarrow$  a | b | c | ...



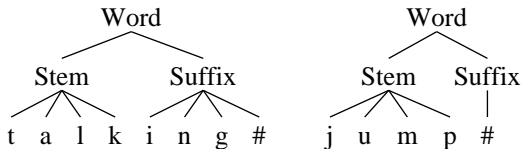
- Grammar's trees can represent any segmentation of words into stems and suffixes

$\Rightarrow$  Can represent true segmentation

- But grammar's *units of generalization* (PCFG rules) are "too small" to learn morphemes

# A “CFG” with one rule per possible morpheme

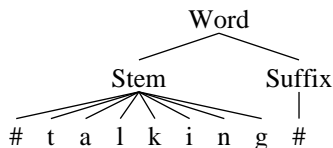
Word  $\rightarrow$  Stem Suffix  
Stem  $\rightarrow$  all possible stems  
Suffix  $\rightarrow$  all possible suffixes



- A rule for each morpheme  
 $\Rightarrow$  “PCFG” can represent probability of each morpheme
- *Unbounded number of possible rules, so this is not a PCFG*
  - ▶ not a practical problem, as only a finite set of rules could possibly be used in any particular data set

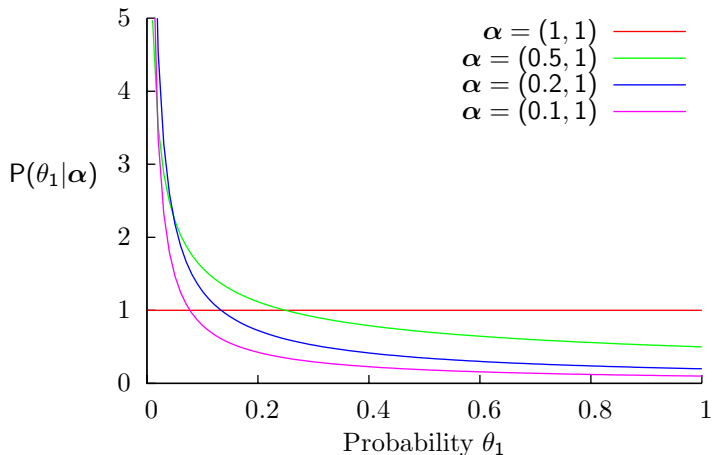
## Maximum likelihood estimate for $\theta$ is trivial

- Maximum likelihood selects  $\theta$  that minimizes KL-divergence between model and training data  $\mathbf{W}$  distributions
  - *Saturated model* in which each word is generated by its own rule replicates training data distribution  $\mathbf{W}$  exactly
- ⇒ Saturated model is maximum likelihood estimate
- Maximum likelihood estimate does not find any suffixes



# Forcing generalization via sparse Dirichlet priors

- Idea: use Bayesian prior that prefers fewer rules
- Set of rules is fixed a priori in Bayesian PCFGs, but can “turn rule off” by setting  $\theta_{A \rightarrow \beta} \approx 0$
- Dirichlet prior with  $\alpha_{A \rightarrow \beta} \approx 0$  prefers  $\theta_{A \rightarrow \beta} \approx 0$



# Morphological segmentation experiment

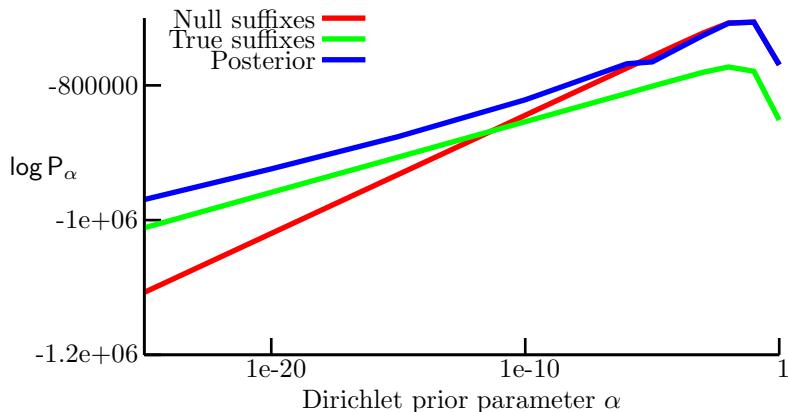
- Trained on orthographic verbs from U Penn. Wall Street Journal treebank
- Uniform Dirichlet prior prefers sparse solutions as  $\alpha \rightarrow 0$
- Metropolis-within-Gibbs sampler used to sample from posterior distribution of parses
  - ▶ reanalyses each word based on a grammar estimated from the parses of the other words



# Posterior samples from WSJ verb tokens

$\alpha = 0.1$	$\alpha = 10^{-5}$	$\alpha = 10^{-10}$	$\alpha = 10^{-15}$
expect	expect	expect	expect
expects	expects	expects	expects
expected	expected	expected	expected
expecting	expect ing	expect ing	expect ing
include	include	include	include
includes	includes	includ es	includ es
included	included	includ ed	includ ed
including	including	including	including
add	add	add	add
adds	adds	adds	add s
added	added	add ed	added
adding	adding	add ing	add ing
continue	continue	continue	continue
continues	continues	continue s	continue s
continued	continued	continu ed	continu ed
continuing	continuing	continu ing	continu ing
report	report	report	report

# Log posterior of models on token data



- Correct solution is nowhere near as likely as posterior  
 $\Rightarrow$  model is wrong!

# Independence assumptions in PCFGs

- Context-free grammars are “context-free” because the possible expansions of each node do not depend on expansions of other nodes
- Probabilistic CFGs extend this by requiring each node expansion to be *statistically independent* (conditioned on the node’s label)
- This is a very strong assumption, which is often false!
- Morphology grammar contains rule:

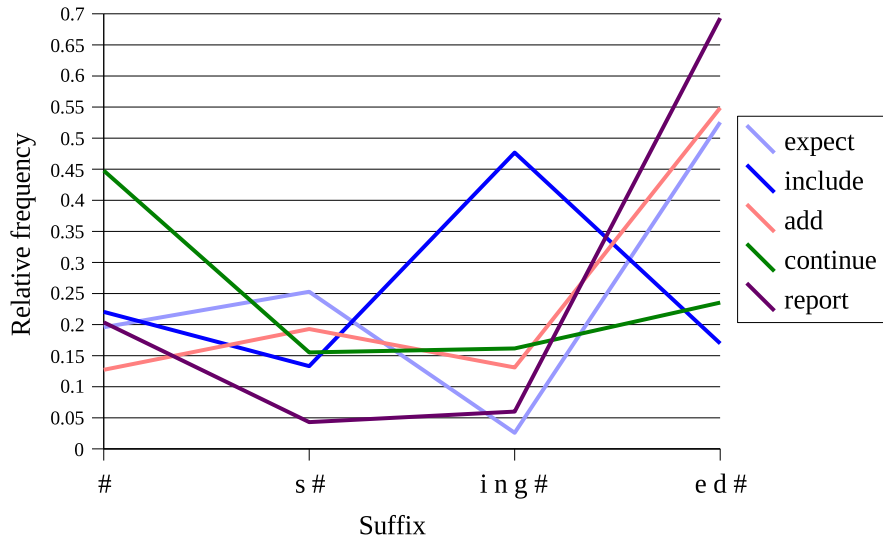
Word  $\rightarrow$  Stem Suffix

- Corresponding independence assumption:

$$P(\text{Word}) = P(\text{Stem}) P(\text{Suffix})$$

causes PCFG model of morphology to fail

# Relative frequencies of inflected verb forms



# Types and tokens

- A word *type* is a distinct word shape
- A word *token* is an occurrence of a word

Data = “the cat chased the other cat”

Tokens = “the”, “cat”, “chased”, “the”, “other”, “cat”

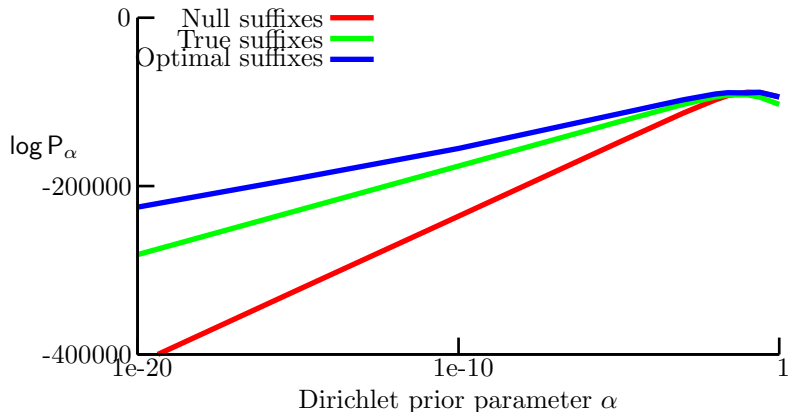
Types = “the”, “cat”, “chased”, “other”

- Estimating  $\theta$  from *word types* rather than word tokens eliminates (most) frequency variation
  - ▶ 4 common verb suffixes, so when estimating from verb types  $\theta_{\text{Suffix} \rightarrow \text{ing}\#} \approx 0.25$
- Several psycholinguists believe that humans learn morphology from word types
- Goldwater et al investigated a morphology-learning model that learnt from an interpolation of types and tokens

# Posterior samples from WSJ verb *types*

$\alpha = 0.1$	$\alpha = 10^{-5}$	$\alpha = 10^{-10}$	$\alpha = 10^{-15}$
expect	expect	expect	exp ect
expects	expect s	expect s	exp ects
expected	expect ed	expect ed	exp ected
expect ing	expect ing	expect ing	exp ecting
include	includ e	includ e	includ e
include s	includ es	includ es	includ es
included	includ ed	includ ed	includ ed
including	includ ing	includ ing	includ ing
add	add	add	add
adds	add s	add s	add s
add ed	add ed	add ed	add ed
adding	add ing	add ing	add ing
continue	continu e	continu e	continu e
continue s	continu es	continu es	continu es
continu ed	continu ed	continu ed	continu ed
continuing	continu ing	continu ing	continu ing
report	report	repo rt	rep ort

# Log posterior of models on type data



- Correct solution is close to optimal at  $\alpha = 10^{-3}$

# Desiderata for an extension of PCFGs

- PCFG rules are “too small” to be effective units of generalization
  - ⇒ generalize over groups of rules
  - ⇒ units of generalization should be chosen based on data
- Type-based inference mitigates non-context-free dependencies
  - ⇒ Hierarchical Bayesian model where:
    - ▶ context-free rules generate types
    - ▶ another process replicates types to produce tokens
- *Adaptor grammars*:
  - ▶ learn probability of entire subtrees (how a nonterminal expands to terminals)
  - ▶ use grammatical hierarchy to define a Bayesian hierarchy, from which type-based inference emerges



# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

# Dirichlet-Multinomials with many outcomes

- Dirichlet prior  $\boldsymbol{\alpha}$ , observed data  $\mathbf{X} = (X_1, \dots, X_n)$

$$P(X_{n+1} = k \mid \mathbf{X}, \boldsymbol{\alpha}) \propto \alpha_k + n_k(\mathbf{X})$$

- Consider a sequence of Dirichlet-multinomials where:
  - ▶ total Dirichlet pseudocount is fixed  $\alpha_{\bullet} = \sum_{k=1}^m \alpha_k$ , and
  - ▶ prior uniform over outcomes  $1, \dots, m$ , so  $\alpha_k = \alpha_{\bullet}/m$
  - ▶ number of outcomes  $m \rightarrow \infty$

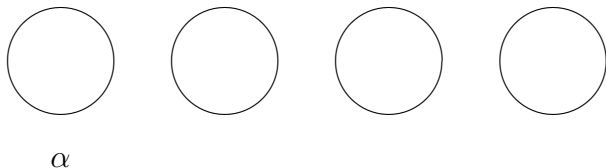
$$P(X_{n+1} = k \mid \mathbf{X}, \alpha_{\bullet}) \propto \begin{cases} n_k(\mathbf{X}) & \text{if } n_k(\mathbf{X}) > 0 \\ \alpha_{\bullet}/m & \text{if } n_k(\mathbf{X}) = 0 \end{cases}$$

But when  $m \gg n$ , most  $k$  are unoccupied (i.e.,  $n_k(\mathbf{X}) = 0$ )

$\Rightarrow$  *Probability of a previously seen outcome  $k \propto n_k(\mathbf{X})$*

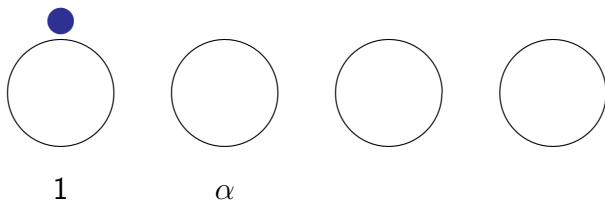
*Probability of an outcome never seen before  $\propto \alpha_{\bullet}$*

# The Chinese restaurant process (1)



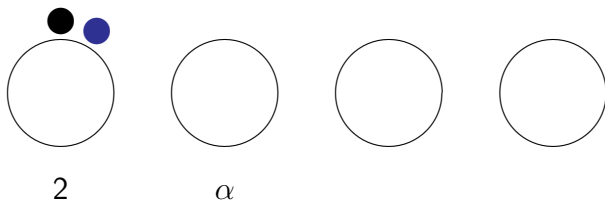
- Each “table” seats infinite number of “customers” (samples)
  - Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
    - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
    - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- $\Rightarrow$  “Rich get richer” power-law dynamics

## The Chinese restaurant process (2)



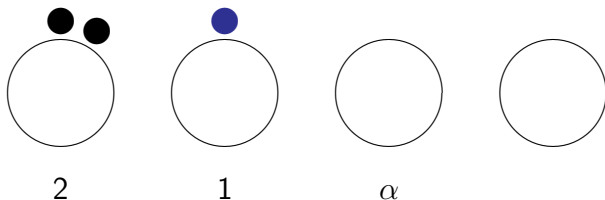
- Each “table” seats infinite number of “customers” (samples)
  - Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
    - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
    - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- $\Rightarrow$  “Rich get richer” power-law dynamics

## The Chinese restaurant process (3)



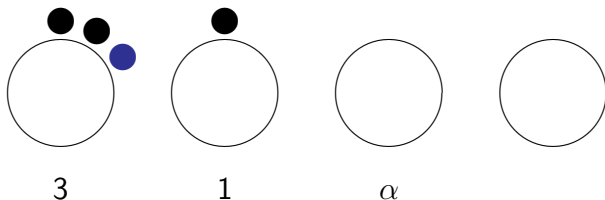
- Each “table” seats infinite number of “customers” (samples)
  - Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
    - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
    - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- $\Rightarrow$  “Rich get richer” power-law dynamics

## The Chinese restaurant process (4)



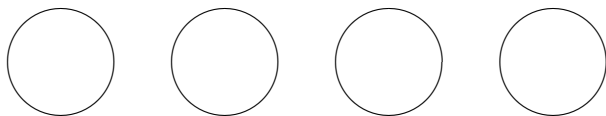
- Each “table” seats infinite number of “customers” (samples)
  - Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
    - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
    - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- $\Rightarrow$  “Rich get richer” power-law dynamics

## The Chinese restaurant process (5)



- Each “table” seats infinite number of “customers” (samples)
  - Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
    - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
    - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- $\Rightarrow$  “Rich get richer” power-law dynamics

## Labeled Chinese restaurant processes (1a)



$\alpha$

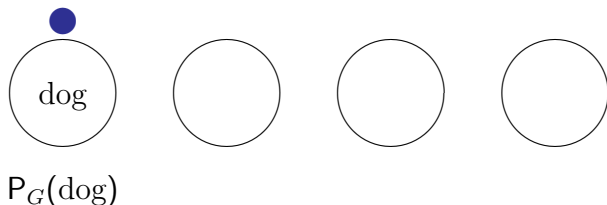
*Generated sequence:*

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*



## Labeled Chinese restaurant processes (1b)

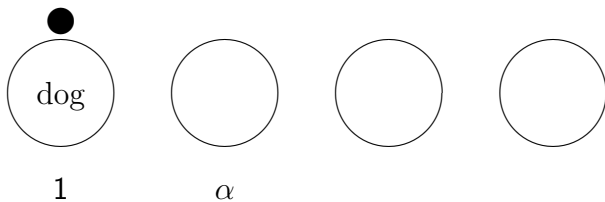


*Generated sequence:* dog

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*

## Labeled Chinese restaurant processes (2a)

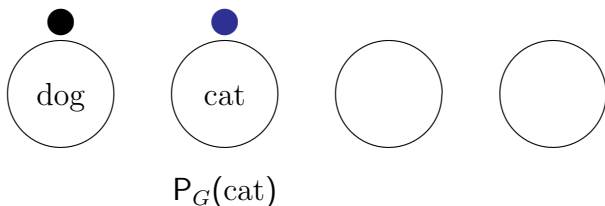


*Generated sequence:* dog

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*

## Labeled Chinese restaurant processes (2b)

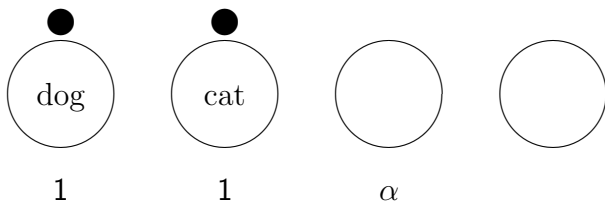


*Generated sequence:* dog, cat

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*

## Labeled Chinese restaurant processes (3a)

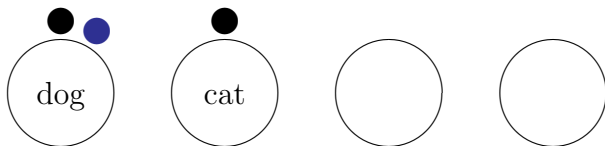


*Generated sequence:* dog, cat

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*

## Labeled Chinese restaurant processes (3b)



*Generated sequence:* dog, cat, dog

- Each occupied table has a label (a “dish”)
- The labels  $Y$  are sampled from  $P_G(Y)$
- Customer  $n + 1$  enters with tables  $1, \dots, m$  occupied:
  - ▶ sits at old table  $k \leq m$  with probability  $\propto n_k$
  - ▶ sits at new table  $k = m + 1$  with probability  $\propto \alpha$
- Emit label (dish) on table
  - ▶ if table doesn't have a label, generate one from  $P_G(Y)$

$\Rightarrow$  only pay probability cost for label *once per table*

# From Chinese restaurants to Dirichlet processes

- Chinese restaurant processes map a distribution  $P_G$  to a stream of samples from a different distribution with the same support
- CRPs specify the *conditional distribution* of the next outcome given the previous ones
- Each CRP run can produce a different distribution over labels
- It defines a mapping from  $\alpha$  and  $P_G$  to a *distribution over distributions*  $DP(\alpha, P_G)$
- $DP(\alpha, P_G)$  is called a *Dirichlet process* (DP) with *concentration parameter*  $\alpha$  and *base distribution*  $P_G$
- The base distribution  $P_G$  can itself be a DP  $\Rightarrow$  *hierarchy* of DPs

# Nonparametric extensions of PCFGs

- Chinese restaurant processes are a nonparametric extension of Dirichlet-multinomials because the number of states (occupied tables) depends on the data
- Two obvious nonparametric extensions of PCFGs:
  - ▶ let the number of nonterminals grow unboundedly
    - refine the nonterminals of an original grammar  
e.g.,  $S_{35} \rightarrow NP_{27} VP_{17}$
    - $\Rightarrow$  infinite PCFG
  - ▶ let the number of rules grow unboundedly
    - “new” rules are compositions of several rules from original grammar
    - equivalent to caching tree fragments
    - $\Rightarrow$  adaptor grammars
- No reason both can't be done together ...

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion



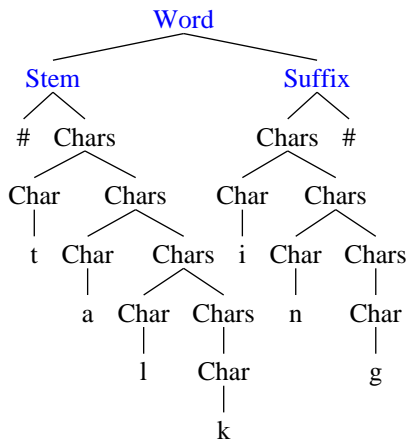
# Adaptor grammars: informal description

- An adaptor grammar has a set of CFG rules
- These determine the possible structures 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)
- Each adapted subtree behaves like a new rule added to the grammar
- The CFG rules of the adapted nonterminals determine the *base distribution* over these trees

# 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*)
- An *unadapted nonterminal*  $A$  expands using  $A \rightarrow \beta$  with probability  $\theta_{A \rightarrow \beta}$
- 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  $\alpha_A \theta_{A \rightarrow \beta}$

# Adaptor grammar morphology example



Word → Stem Suffix

Stem → # Chars

Suffix → #

Suffix → Chars #

Chars → Char

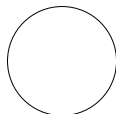
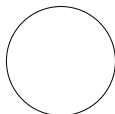
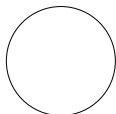
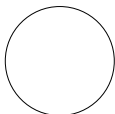
Chars → Char Chars

Char → a | ... | z

- Stem and Suffix rules generate all possible stems and suffixes
- Adapt Word, Stem and Suffix nonterminals
- One Chinese Restaurant process per adapted nonterminal

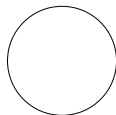
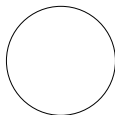
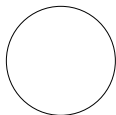
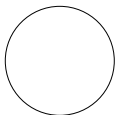
# Morphology adaptor grammar (0)

**Word** restaurant  
Word  $\rightarrow$  Stem Suffix



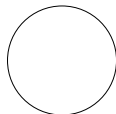
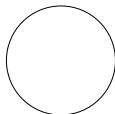
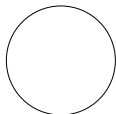
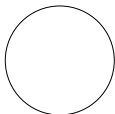
...

**Stem** restaurant  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix** restaurant  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

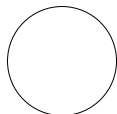
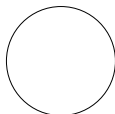
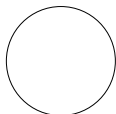
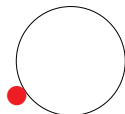


...

**Chars** factory  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a . . . z

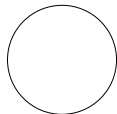
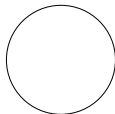
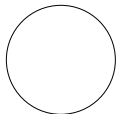
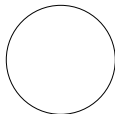
# Morphology adaptor grammar (1a)

**Word** restaurant  
Word  $\rightarrow$  Stem Suffix



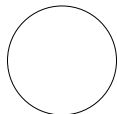
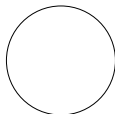
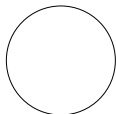
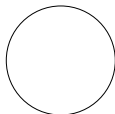
...

**Stem** restaurant  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix** restaurant  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

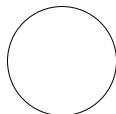
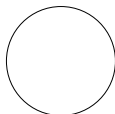
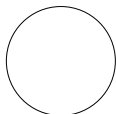
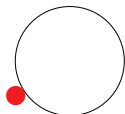


...

**Chars** factory  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a . . . z

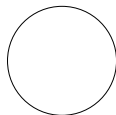
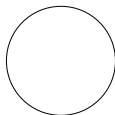
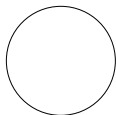
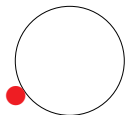
# Morphology adaptor grammar (1b)

**Word** restaurant  
Word  $\rightarrow$  Stem Suffix



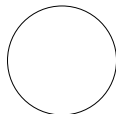
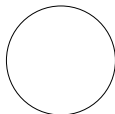
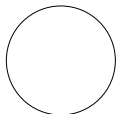
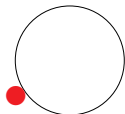
...

**Stem** restaurant  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix** restaurant  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

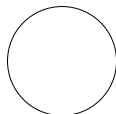
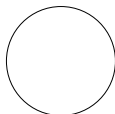
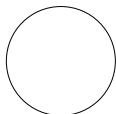
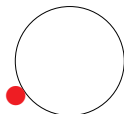


...

**Chars** factory  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a . . . z

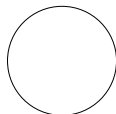
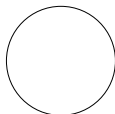
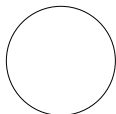
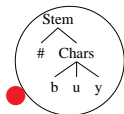
# Morphology adaptor grammar (1c)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



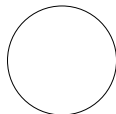
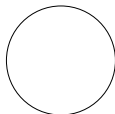
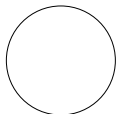
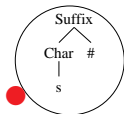
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

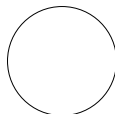
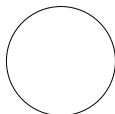
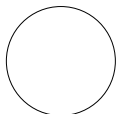
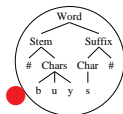


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a . . . z

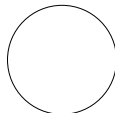
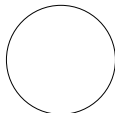
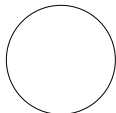
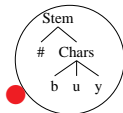
# Morphology adaptor grammar (1d)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



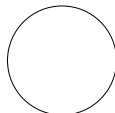
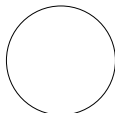
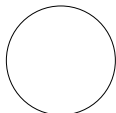
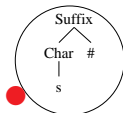
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #



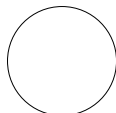
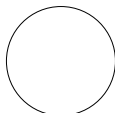
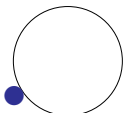
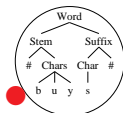
...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z



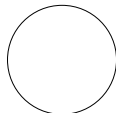
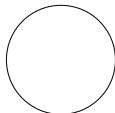
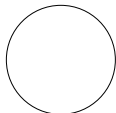
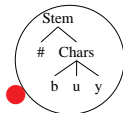
# Morphology adaptor grammar (2a)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



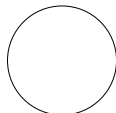
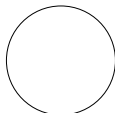
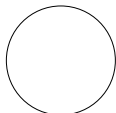
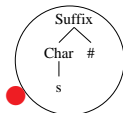
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

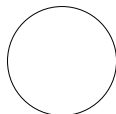
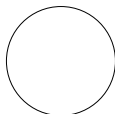
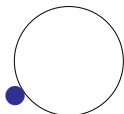
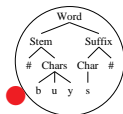


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

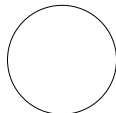
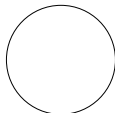
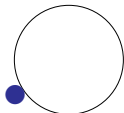
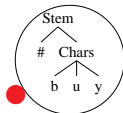
# Morphology adaptor grammar (2b)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



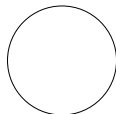
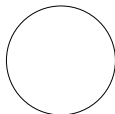
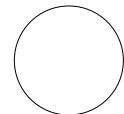
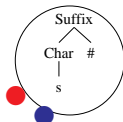
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

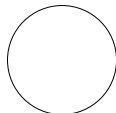
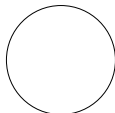
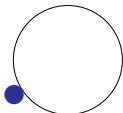
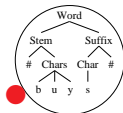


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

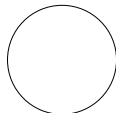
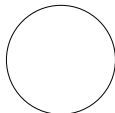
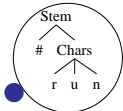
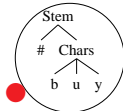
# Morphology adaptor grammar (2c)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



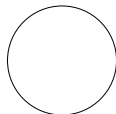
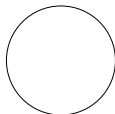
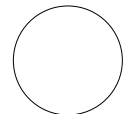
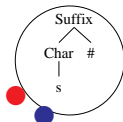
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

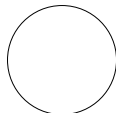
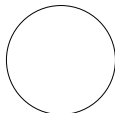
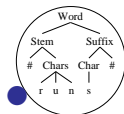
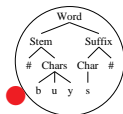


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a . . . z

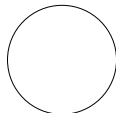
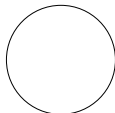
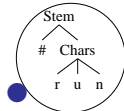
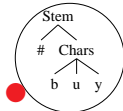
# Morphology adaptor grammar (2d)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



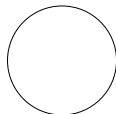
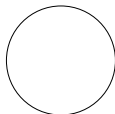
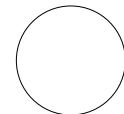
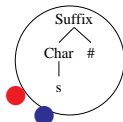
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

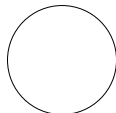
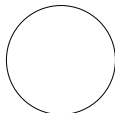
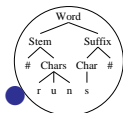
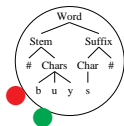


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

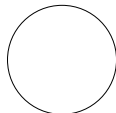
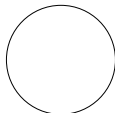
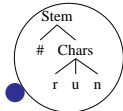
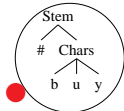
# Morphology adaptor grammar (3)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



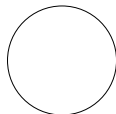
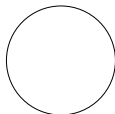
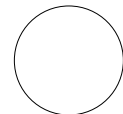
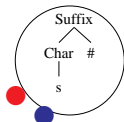
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

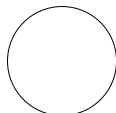
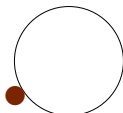
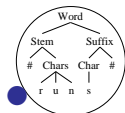
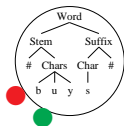


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

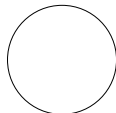
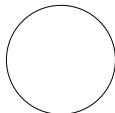
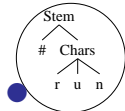
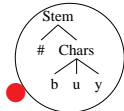
# Morphology adaptor grammar (4a)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



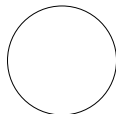
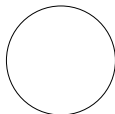
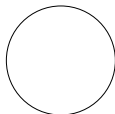
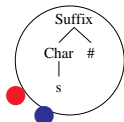
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

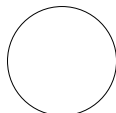
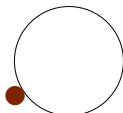
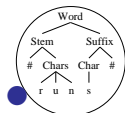
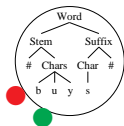


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

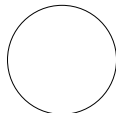
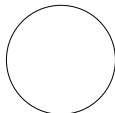
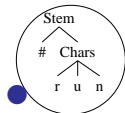
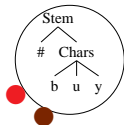
# Morphology adaptor grammar (4b)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



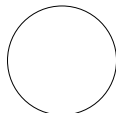
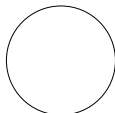
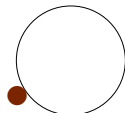
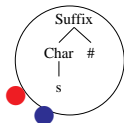
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #

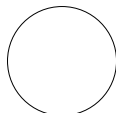
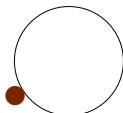
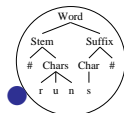
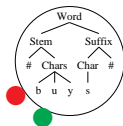


...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

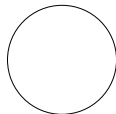
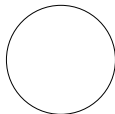
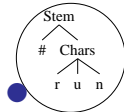
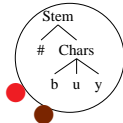
# Morphology adaptor grammar (4c)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



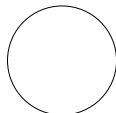
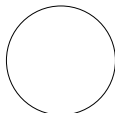
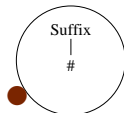
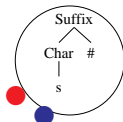
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #



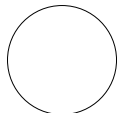
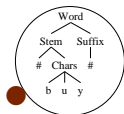
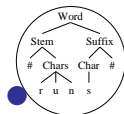
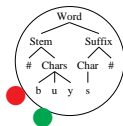
...

**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z



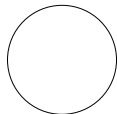
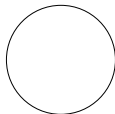
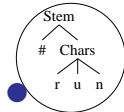
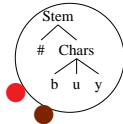
# Morphology adaptor grammar (4d)

**Word restaurant**  
Word  $\rightarrow$  Stem Suffix



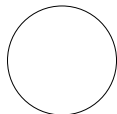
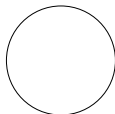
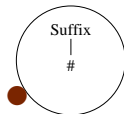
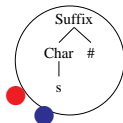
...

**Stem restaurant**  
Stem  $\rightarrow$  #  
Stem  $\rightarrow$  # Chars



...

**Suffix restaurant**  
Suffix  $\rightarrow$  #  
Suffix  $\rightarrow$  Chars #



...

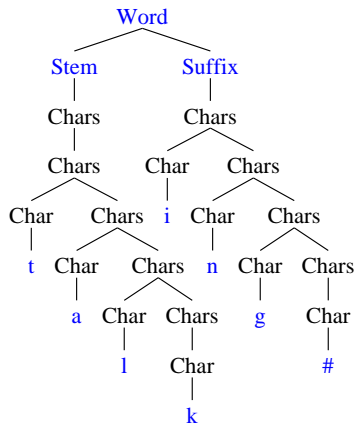
**Chars factory**  
Chars  $\rightarrow$  Char  
Chars  $\rightarrow$  Char Chars  
Char  $\rightarrow$  a...z

# Properties of adaptor grammars

- Possible trees generated by CFG rules  
but the probability of each adapted tree is estimated separately
  - Probability of a subtree  $\tau$  is proportional to:
    - ▶ the number of times  $\tau$  was seen before  
 $\Rightarrow$  “rich get richer” dynamics (Zipf distributions)
    - ▶ plus  $\alpha_A$  times prob. of generating it via PCFG expansion
- $\Rightarrow$  Useful compound structures can be *more probable than their parts*
- PCFG rule probabilities estimated *from table labels*
    - $\Rightarrow$  learns from types, not tokens
    - $\Rightarrow$  dampens frequency variation

# Bayesian hierarchy inverts grammatical hierarchy

- Grammatically, a Word is composed of a Stem and a Suffix, which are composed of Chars
- To generate a new Word from an adaptor grammar
  - ▶ reuse an old Word, or
  - ▶ generate a fresh one from the base distribution, i.e., generate a Stem and a Suffix
- Lower in the tree  
⇒ higher in Bayesian hierarchy



## PCFGs as recursive mixtures

For simplicity assume all rules are of the form  $A \rightarrow BC$  or  $A \rightarrow w$ , where  $A, B, C \in N$  (nonterminals) and  $w \in T$  (terminals).

Each nonterminal  $A \in N$  generates a distribution  $G_A$  over the trees rooted in  $A$ .

$$G_A = \sum_{A \rightarrow BC \in R_A} \theta_{A \rightarrow BC} \text{TREE}_A(G_B, G_C) + \sum_{A \rightarrow w \in R_A} \theta_{A \rightarrow w} \text{TREE}_A(w)$$

where  $\text{TREE}_A(w)$  puts all of its mass on the tree with child  $w$  and  $\text{TREE}_A(P, Q)$  is the distribution over trees rooted in  $A$  with children distributed according to  $P$  and  $Q$  respectively.

$$\text{TREE}_A(P, Q) \left( \begin{array}{c} A \\ \swarrow \quad \searrow \\ t_1 \quad t_2 \end{array} \right) = P(t_1) Q(t_2)$$

The tree language generated by the PCFG is  $G_S$ .

## Adaptor grammars as recursive mixtures

An adaptor grammar  $(G, \boldsymbol{\theta}, \boldsymbol{\alpha})$  is a PCFG  $(G, \boldsymbol{\theta})$  together with a parameter vector  $\boldsymbol{\alpha}$  where for each  $A \in N$ ,  $\alpha_A$  is the parameter of the Dirichlet process associated with  $A$ .

$$\begin{aligned} G_A &\sim \text{DP}(\alpha_A, H_A) \text{ if } \alpha_A > 0 \\ &= H_A \quad \text{if } \alpha_A = 0 \end{aligned}$$

$$H_A = \sum_{A \rightarrow BC \in R_A} \theta_{A \rightarrow BC} \text{TREE}_A(G_B, G_C) + \sum_{A \rightarrow w \in R_A} \theta_{A \rightarrow w} \text{TREE}_A(w)$$

The grammar generates the distribution  $G_S$ .

There is one Dirichlet Process for each non-terminal  $A$  where  $\alpha_A > 0$ . Its base distribution  $H_A$  is a mixture of the language generated by the Dirichlet processes associated with other non-terminals.

# Bayesian priors on adaptor grammar parameters

- Parameters of adaptor grammars:
  - ▶ probabilities  $\theta_{A \rightarrow \beta}$  of base grammar rules  $A \rightarrow \beta$
  - ▶ concentration parameters  $\alpha_A$  of adapted nonterminals  $A$
- Put Bayesian priors on these parameters
  - ▶ (Uniform) Dirichlet prior on base grammar rule probabilities  $\theta$
  - ▶ Vague Gamma prior on concentration parameter on  $\alpha_A$
- We also use a generalization of CRPs called “Pitman-Yor processes”, and put a uniform Dirichlet prior on its  $a$  parameter

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

**Word segmentation using adaptor grammars**

Bayesian inference of adaptor grammars

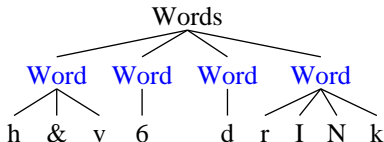
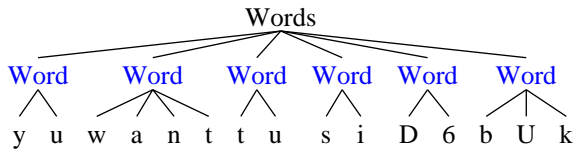
Conclusion

# Unigram model of word segmentation

- Unigram model: each word is generated independently
- Input is *unsegmented broad phonemic transcription* (Brent)  
Example: y u w a n t t u s i D 6 b u k
- Adaptor for **Word** non-terminal caches previously seen words

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

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

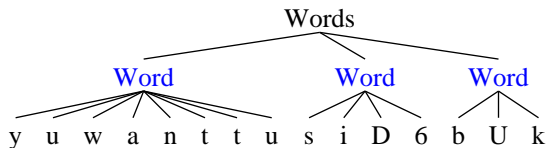
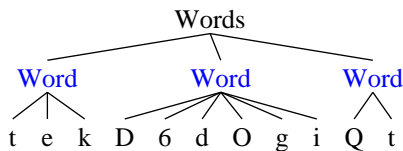


- Unigram word segmentation on Brent corpus: 55% token f-score



# Unigram model often finds collocations

- 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



# Unigram word segmentation grammar learnt

- Based on the base grammar rules

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

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

the adapted grammar contains 1,712 rules such as:

15758 Words  $\rightarrow$  Word Words

9791 Words  $\rightarrow$  Word

1660 Word  $\rightarrow$  Phoneme<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

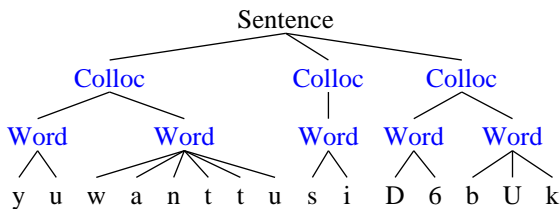
20 Word  $\rightarrow$  l n D 6 h Q s

# Modeling collocations improves segmentation

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

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

Word  $\rightarrow$  Phoneme<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 (75% token f-score; same as Goldwater's bigram model)
- Two levels of Collocations improves slightly (76%)

# Syllables + Collocations + Word segmentation

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

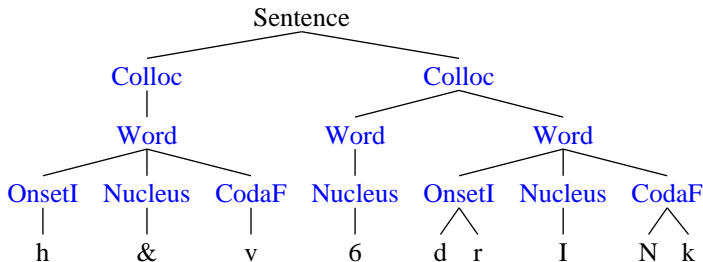
Word  $\rightarrow$  SyllableIF

Word  $\rightarrow$  SyllableI Syllable SyllableF

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

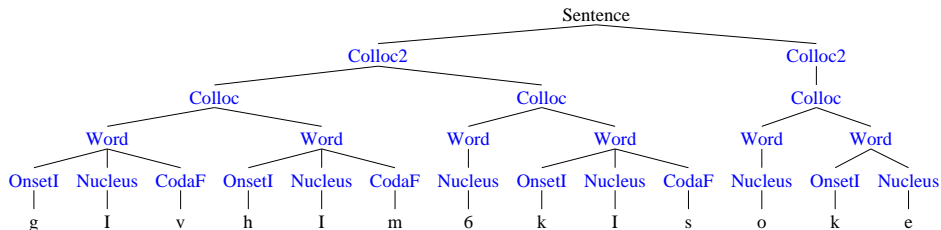
Word  $\rightarrow$  SyllableI SyllableF

Syllable  $\rightarrow$  (as before)



- Word segmentation f-score = 68%
- With 2 Collocation levels f-score = 84% (better than Colloc+morphology)
- Without distinguishing initial/final clusters f-score = 82%

# Syllables + 2-level Collocations + Word segmentation



## Word segmentation results summary

	Collocation levels			
	0	1	2	3
word	0.55	0.73	0.75	0.72
morphology	0.35	0.55	0.79	0.73
syllable	0.32	0.69	0.82	0.79
syllableIF	0.46	0.68	<b>0.84</b>	<b>0.84</b>

- We can learn collocations and syllable structure together with word segmentation, even though we don't know where the word boundaries are
- Learning these together improves word segmentation accuracy
  - ▶ are there other examples of *synergistic interaction* in language learning?

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

# Estimating adaptor grammars

- Need to estimate:
  - ▶ table labels (subtrees)  $\tau$  and customer count for each table
  - ▶ (optional) probabilities of base grammar productions
  - ▶ (optional) DP parameters  $\alpha$
- Component-wise Metropolis-within-Gibbs sampler
  - ▶  $i$ th component is the parse tree  $T_i$  for input string  $W_i$
  - ▶ sample parse  $T_i$  for  $W_i$  using grammar  $G'(\mathbf{T}_{-i})$  estimated from parses  $\mathbf{T}_{-i}$  for other inputs
- Sampling directly from conditional distribution of parses seems intractable
  - ▶ construct PCFG proposal grammar  $G'(\mathbf{T}_{-i})$  on the fly
  - ▶ each table label  $\tau$  corresponds to a production in PCFG approximation
  - ▶ Use accept/reject to convert samples from PCFG approx to samples from adaptor grammar



# PCFG proposal grammar

- Recall that in a CRP,
  - ▶ picking existing table  $\tau$  with prob.  $\propto n_\tau$  (number of customers seated at  $\tau$ )
  - ▶ picking new table with prob.  $\propto \alpha$  (DP concentration parameter)
- Rules of PCFG proposal grammar  $G'$  consist of:
  - ▶ rules  $A \rightarrow \beta$  from base PCFG:  $\theta'_{A \rightarrow \beta} \propto \alpha_A \theta_{A \rightarrow \beta}$
  - ▶ A rule  $A \rightarrow \text{YIELD}(\tau)$  for each table  $\tau$  in  $A$ 's restaurant:  
 $\theta'_{A \rightarrow \text{YIELD}(\tau)} \propto n_\tau$ , the number of customers at table  $\tau$
- Parses of  $G'$  can be mapped back to adaptor grammar parses

# Implementation details

- The sampler just described can take a long time to “burn in”
- It’s hard to predict when annealing will help
- Resampling the table labels (after each pass through the sentences) dramatically increases mobility
  - ▶ makes it possible to reanalyse many sentences at once

# Outline

Bayesian inference of multinomials

Inference of Probabilistic Context-Free Grammars

Limitations of PCFG learning

Chinese restaurant processes and Dirichlet processes

Adaptor grammars

Word segmentation using adaptor grammars

Bayesian inference of adaptor grammars

Conclusion

## Summary and future work

- Adaptor grammars “adapt” their distribution to the strings they have generated
- They learn the probabilities of the subtrees of the adapted nonterminals
- This makes adaptor grammars *non-parametric*; the number of subtrees they track depends on the data
- A variety of different linguistic phenomena can be described with adaptor grammars
- Because they are grammars, they are easy to design and compose
- But they still have a “context-freeness” that makes it impossible to express e.g., Goldwater’s bigram word segmentation model. Can we add context-sensitivity in a manageable way?
- The MCMC sampling algorithm used does not seem to scale well to large data or complicated grammars. Are there better estimators?