

# Parsing and Learning

International Workshop on Parsing Technology

Mark Johnson

Brown University / Macquarie University

October 2009

# Parsing and Learning

- Parsing and grammar induction traditionally viewed as distinct
  - ▶ Unknown words and phrases  $\Rightarrow$  parser must be prepared to adapt
- Most grammar learning systems are *error-driven*
  - ▶ parsing is a central component of most learning algorithms
- Traditional statistical learning algorithms learn *rule probabilities*
- Learning grammar rules
  - ▶ via “generate and test”
  - ▶ via *non-parametric Bayes*  $\Rightarrow$  *adaptor grammars*

# Talk outline

- Learning PCFG rule probabilities from terminal strings
  - ▶ collapsed Gibbs samplers explicitly sample parses, but don't explicitly represent rule probabilities
- Non-parametric extensions to PCFGs with unboundedly many possible rules
  - ▶ adaptor grammars learn probabilities of arbitrary subtrees
  - ▶ collapsed Gibbs sampler only needs to represent sample parses, and not probabilities of infinitely many rules
- Adaptor grammars can be used to learn:
  - ▶ simple concatenative morphology
  - ▶ unsupervised word segmentation
  - ▶ collocations in topic models
  - ▶ structure in names

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# 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

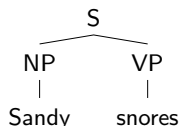
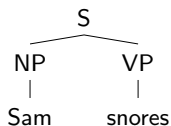
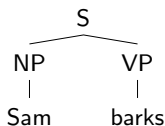
Rule $r$	$\theta_r$
$NP \rightarrow \text{Sandy}$	0.25
$VP \rightarrow \text{snores}$	0.4

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

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

# Maximum likelihood estimation from visible parses

- Each rule expansion is sampled from parent's multinomial
- ⇒ Maximum Likelihood Estimator (MLE) is rule's *relative frequency*



Rule $r$	$n_r$	$\theta_r$
$S \rightarrow NP VP$	3	1.0
$NP \rightarrow Sam$	2	0.66
$VP \rightarrow barks$	1	0.33

Rule $r$	$n_r$	$\theta_r$
$NP \rightarrow Sandy$	1	0.33
$VP \rightarrow snores$	2	0.66

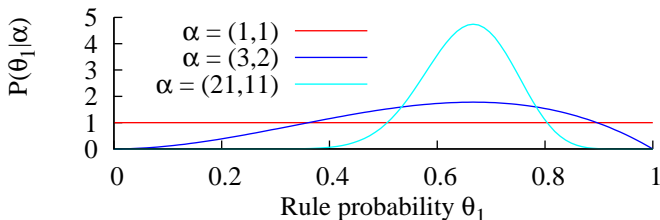
- But MLE is often *overly certain*, especially with sparse data
  - ▶ E.g., “accidental zeros”  $n_r = 0 \Rightarrow \theta_r = 0$ .

# Bayesian estimation from visible parses

- Bayesian estimators estimate a *distribution* over rule probabilities

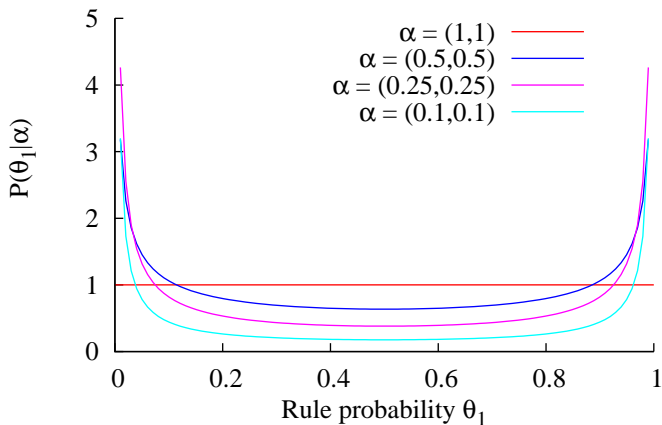
$$\underbrace{P(\boldsymbol{\theta} | \mathbf{n})}_{\text{Posterior}} \propto \underbrace{P(\mathbf{n} | \boldsymbol{\theta})}_{\text{Likelihood}} \underbrace{P(\boldsymbol{\theta})}_{\text{Prior}}$$

- Dirichlet distributions* are conjugate priors for multinomials
  - A Dirichlet distribution over  $(\theta_1, \dots, \theta_m)$  is specified by positive parameters  $(\alpha_1, \dots, \alpha_m)$
  - If Prior =  $\text{Dir}(\boldsymbol{\alpha})$  then Posterior =  $\text{Dir}(\boldsymbol{\alpha} + \mathbf{n})$



# Sparse Dirichlet priors

- As  $\alpha \rightarrow 0$ , Dirichlet distributions become peaked around 0  
“Grammar includes some of these rules, but we don't know which!”





# Estimating rule probabilities from strings alone

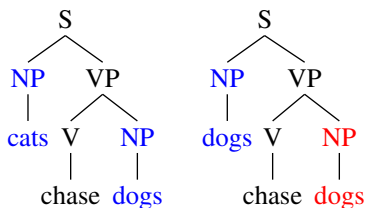
- Input: *terminal strings* and *grammar rules*
- Output: *rule probabilities*  $\theta$
- In general, no closed-form solution for  $\theta$ 
  - ▶ iterative algorithms usually involving *repeatedly reparsing training data*
- *Expectation Maximization* (EM) procedure generalizes visible data ML estimators to hidden data problems
- *Inside-Outside algorithm* is a cubic-time EM algorithm for PCFGs
- Bayesian estimation of  $\theta$  via:
  - ▶ *Variational Bayes* or
  - ▶ *Markov Chain Monte Carlo* (MCMC) methods such as *Gibbs sampling*

# Gibbs sampler for parse trees and rule probabilities

- Input: *terminal strings*  $(x_1, \dots, x_n)$ , *grammar rules* and *Dirichlet prior parameters*  $\alpha$
- Output: stream of sample *rule probabilities*  $\theta$  and *parse trees*  $t = (t_1, \dots, t_n)$
- Algorithm:
  - Assign parse trees to the strings somehow (e.g., randomly)
  - Repeat forever:
    - Compute rule counts  $n$  from  $t$
    - Sample  $\theta$  from  $\text{Dir}(\alpha + n)$
    - For each string  $x_i$ :
      - replace  $t_i$  with a tree sampled from  $P(t|x_i, \theta)$ .
- After *burn-in*,  $(\theta, t)$  are distributed according to Bayesian posterior
- Sampling parse tree from  $P(t|x_i, \theta)$  involves parsing string  $x_i$ .

# Collapsed Gibbs samplers

- *Integrate out* rule probabilities  $\theta$  to obtain *predictive distribution*  $P(t_i|x_i, t_{-i})$  of parse  $t_i$  for sentence  $x_i$  given other parses  $t_{-i}$
- *Collapsed Gibbs sampler*
  - For each sentence  $x_i$  in training data:
    - Replace  $t_i$  with a sample from  $P(t|x_i, t_{-i})$
- A problem:  $P(t_i|x_i, t_{-i})$  is *not a PCFG distribution*
  - $\Rightarrow$  no dynamic-programming sampler (AFAIK)



# Metropolis-Hastings samplers

- Metropolis-Hastings (MH) acceptance-rejection procedure uses samples from a *proposal distribution* to produce samples from a *target distribution*
- When sentence size  $\ll$  training data size,  $P(t_i|x_i, t_{-i})$  is almost a PCFG distribution
  - ▶ use a PCFG approximation based on  $t_{-i}$  as *proposal distribution*
  - ▶ apply MH to transform proposals to  $P(t_i|x_i, t_{-i})$
- To construct a Metropolis-Hastings sampler you need to be able to:
  - ▶ efficiently sample from proposal distribution
  - ▶ calculate ratios of parse probabilities under proposal distribution
  - ▶ calculate ratios of parse probabilities under target distribution

# Collapsed Metropolis-within-Gibbs sampler for PCFGs

- Input: *terminal strings*  $(x_1, \dots, x_n)$ , *grammar rules* and *Dirichlet prior parameters*  $\alpha$
- Output: stream of sample *parse trees*  $t = (t_1, \dots, t_n)$
- Algorithm:
  - Assign parse trees to the strings somehow (e.g., randomly)
  - Repeat forever:
    - For each sentence  $x_i$  in training data:
      - Compute rule counts  $n_{-i}$  from  $t_{-i}$
      - Compute proposal grammar probabilities  $\theta$  from  $n_{-i}$
      - Sample a tree  $t$  from  $P(t|x_i, \theta)$
      - Replace  $t_i$  with  $t$  according to
        - Metropolis-Hastings accept-reject formula

## Q: Are there efficient *local-move* tree samplers?

- DP PCFG samplers require  $O(n^3)$  time per sample; are there faster samplers?
- For HMMs, a *point-wise* sampler resamples the state labeling *one word* at a time
  - ▶ All state *sequences* reachable by one or more such moves
  - ▶ Transition probabilities don't require dynamic programming  $\Rightarrow$  no need for MH
- Question: *Are there efficiently-computable local moves that can transform any parse into any other parse of same string?*
  - ▶ for HMMs, DP sampler generally more effective than point-wise sampler
  - ▶ point-wise samplers could be more effective for grammars with complex DP parsing algorithms (e.g., TAG, CCG, LFG, HPSG)

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# Learning grammar rules

- Input: *terminal strings*
- Output: *grammar rules* and *rule probabilities  $\theta$*
- “Generate and test” approach (Carroll and Charniak, Stolcke)
  - Guess an initial set of rules
  - Repeat:
    - re-estimate rule probabilities from strings
    - prune low probability rules
    - propose additional potentially useful rules*
- Non-parametric Bayesian methods seem to provide a more systematic approach



# Non-parameteric Bayesian extensions to PCFGs

- Non-parametric  $\Rightarrow$  no fixed set of parameters
- Two obvious non-parametric extensions to PCFGs:
  - ▶ let the set of non-terminals grow unboundedly
    - given an initial grammar with coarse-grained categories, split non-terminals into more refined categories  
 $S_{12} \rightarrow NP_7 VP_4$  instead of  $S \rightarrow NP VP$ .
    - PCFG generalization of “infinite HMM”.
  - ▶ let the set of rules grow unboundedly  $\Rightarrow$  *adaptor grammars*
    - use a (meta-)grammar to generate potential rules
    - learn subtrees and their probabilities  
i.e., tree substitution grammar, where we learn the fragments as well as their probabilities
- No reason both can't be done at once ...

# Learning syntax is hard!

- Can formulate learning syntax as Bayesian estimation
- On toy data, Bayesian estimators do well
- Results are disappointing on “real” data
  - ▶ wrong algorithm?
  - ▶ wrong kind of grammar?
  - ▶ wrong type of training data?
- This paper focuses on learning grammars for simpler phenomena:
  - ▶ Morphological segmentation (e.g., *walking* = *walk+ing*)
  - ▶ Word segmentation of unsegmented phoneme sequences
  - ▶ Learning collocations in topic models
  - ▶ Learning internal structure of named-entity NPs

# A CFG for stem-suffix morphology

Word  $\rightarrow$  Stem Suffix

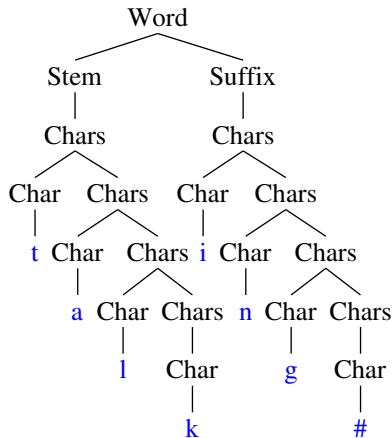
Stem  $\rightarrow$  Chars

Suffix  $\rightarrow$  Chars

Chars  $\rightarrow$  Char

Chars  $\rightarrow$  Char Chars

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



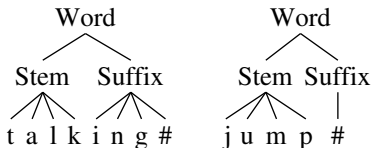
- Grammar's trees can represent any segmentation of words into stems and suffixes

$\Rightarrow$  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

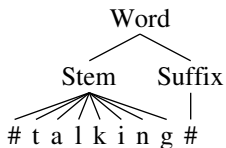
Word → Stem Suffix  
Stem → *all possible stems*  
Suffix → *all possible suffixes*



- A rule for each morpheme  
⇒ “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



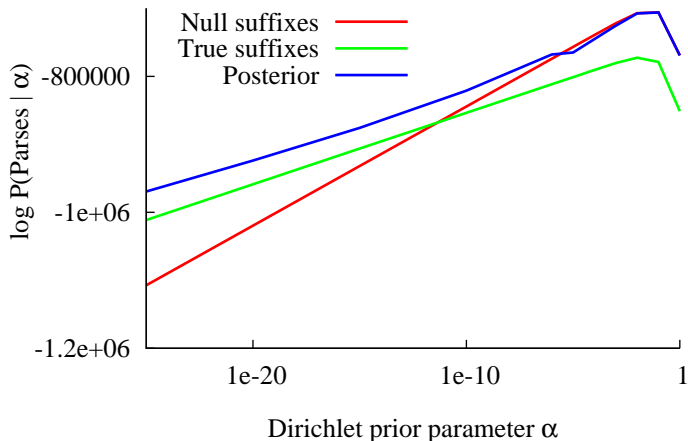
# Forcing generalization using Dirichlet priors

- Maximum Likelihood solution analyses each word as a separate stem
  - ▶ fails to generalize
  - ▶ one non-zero probability rule per word type in data
- Dirichlet prior prefers  $\theta = 0$  when  $\alpha \rightarrow 0$ 
  - ▶ use Dirichlet prior to prefer *sparse rule probability vectors*
- Following experiments use orthographic verbs from U Penn. WSJ treebank

# 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

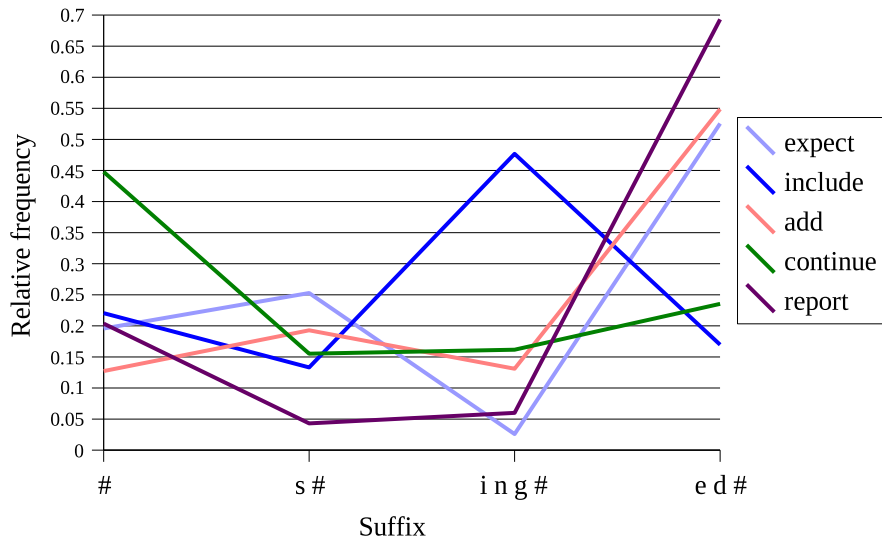
# Log posterior for models on token data



- Correct solution is nowhere near as likely as posterior  
⇒ model is wrong!



# 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”

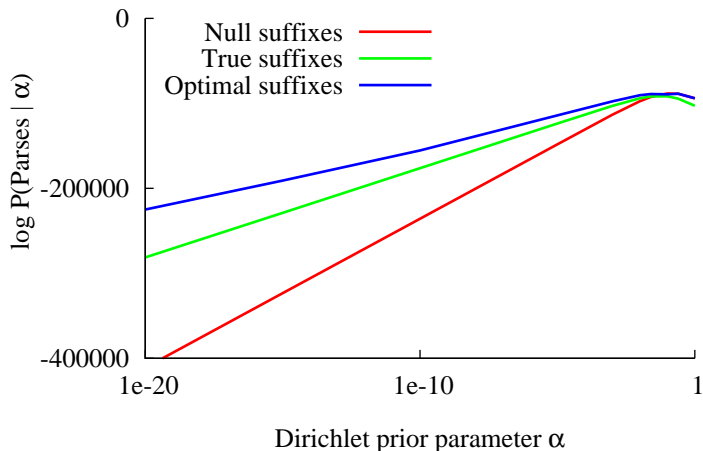
Types = “the”, “cat”, “chased”, “other”

- 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 over-dispersion
  - ⇒ 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
  - ▶ inspired by Sharon Goldwater’s models

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

## Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# From Dirichlet-multinomials to Chinese Restaurant Processes

- Observations  $\mathbf{z} = (z_1, \dots, z_n)$  ranging over outcomes  $1, \dots, m$
- Outcome  $k$  observed  $n_k(\mathbf{z})$  times in data  $\mathbf{z}$
- *Predictive distribution* with uniform Dirichlet prior:

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

- Let  $m \rightarrow \infty$

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

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

- If outcomes are exchangeable  $\Rightarrow$  number in order of occurrence  
 $\Rightarrow$  *Chinese Restaurant Process*

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

# Chinese Restaurant Process (0)

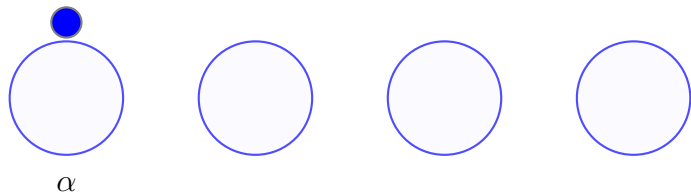


- Customer  $\rightarrow$  table mapping  $z =$
- $P(z) = 1$
- Next customer chooses a table according to:

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



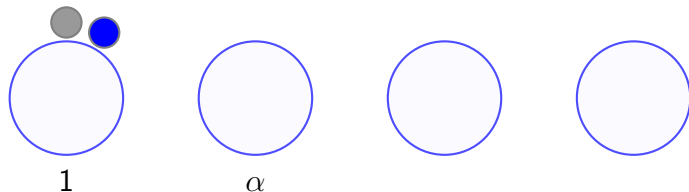
# Chinese Restaurant Process (1)



- Customer  $\rightarrow$  table mapping  $z = 1$
- $P(z) = \alpha/\alpha$
- Next customer chooses a table according to:

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

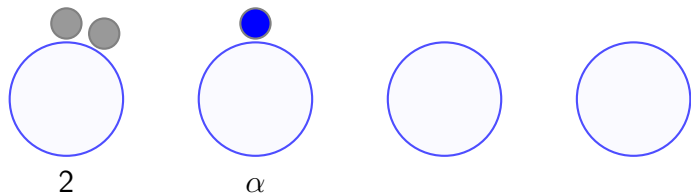
## Chinese Restaurant Process (2)



- Customer  $\rightarrow$  table mapping  $z = 1, 1$
- $P(z) = \alpha/\alpha \times 1/(1 + \alpha)$
- Next customer chooses a table according to:

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

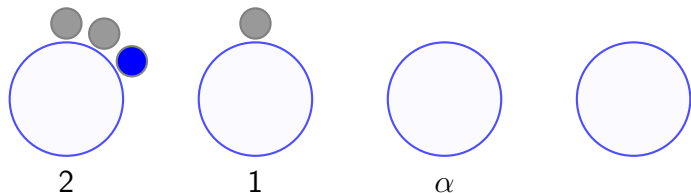
## Chinese Restaurant Process (3)



- 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:

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

## Chinese Restaurant Process (4)



- Customer  $\rightarrow$  table mapping  $\mathbf{z} = 1, 1, 2, 1$
- $P(\mathbf{z}) = \frac{\alpha}{\alpha} \times \frac{1}{(1 + \alpha)} \times \frac{\alpha}{(2 + \alpha)} \times \frac{2}{(3 + \alpha)}$
- Next customer chooses a table according to:

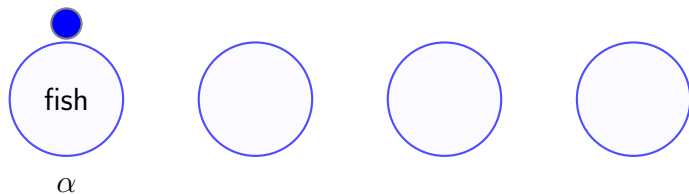
$$P(Z_{n+1} = k \mid \mathbf{z}) \propto \begin{cases} n_k(\mathbf{z}) & \text{if } k \leq m = \max(\mathbf{z}) \\ \alpha & \text{if } k = m + 1 \end{cases}$$

# Labeled Chinese Restaurant Process (0)



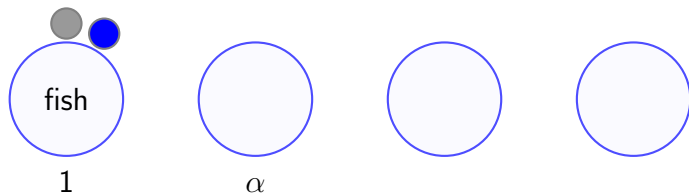
- Table  $\rightarrow$  label mapping  $\mathbf{y} =$
- Customer  $\rightarrow$  table mapping  $\mathbf{z} =$
- Output sequence  $\mathbf{x} =$
- $P(\mathbf{x}) = 1$
  
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$
- All customers sitting at table  $k$  (i.e.,  $z_i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

# Labeled Chinese Restaurant Process (1)



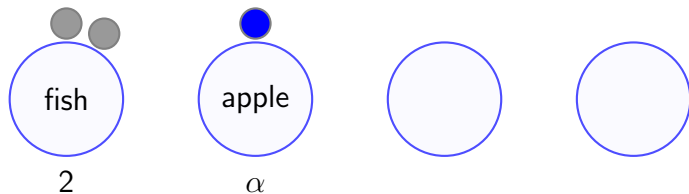
- Table  $\rightarrow$  label mapping  $\mathbf{y} = \text{fish}$
- Customer  $\rightarrow$  table mapping  $z = 1$
- Output sequence  $\mathbf{x} = \text{fish}$
- $P(\mathbf{x}) = \alpha/\alpha \times P_0(\text{fish})$
  
- *Base distribution*  $P_0(Y)$  generates a *label*  $y_k$  for each table  $k$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

## Labeled 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  $x_i = y_{z_i}$

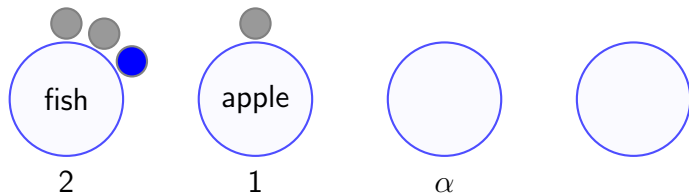
## Labeled Chinese Restaurant Process (3)



- 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  $x_i = y_{z_i}$



## Labeled 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  $x_i = y_{z_i}$

# Summary: Chinese Restaurant 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
- *Labeled CRPs and PYPs* use a *base distribution* to label each table
  - ▶ base distribution can have *infinite support*
  - ▶ concentrates mass on a countable subset

# Labeled Chinese restaurants and Dirichlet processes

- A labeled Chinese restaurant processes maps a *base distribution*  $P_B$  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  $P_B$  to a *distribution over distributions*  $DP(\alpha, P_B)$
- $DP(\alpha, P_B)$  is called a *Dirichlet process* (DP) with *concentration parameter*  $\alpha$  and *base distribution*  $P_B$
- The base distribution  $P_B$  can itself be defined by a DP  
 $\Rightarrow$  *hierarchy* of DPs

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

**Adaptor grammars**

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# 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*
  - ▶ each adapted nonterminal  $A$  has a *concentration parameter*  $\alpha_A$
- 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 for stem-suffix morphology (0)

Word  $\rightarrow$  Stem Suffix



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



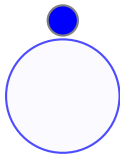
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



Generated words:

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

Word  $\rightarrow$  Stem Suffix



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



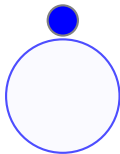
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



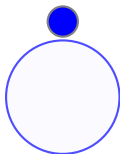
Generated words:

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

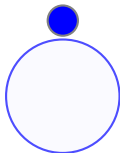
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



Suffix → Phoneme<sup>\*</sup>

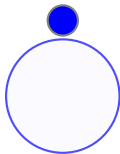


Generated words:

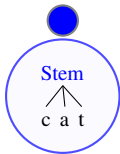


# Adaptor grammar for stem-suffix morphology (1c)

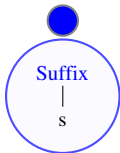
Word  $\rightarrow$  Stem Suffix



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



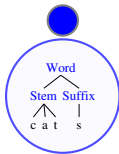
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



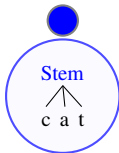
Generated words:

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

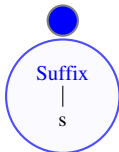
Word  $\rightarrow$  Stem Suffix



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



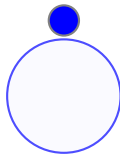
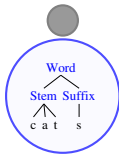
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



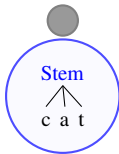
Generated words: **cats**

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

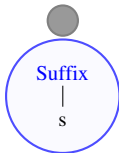
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



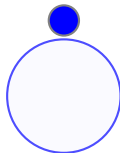
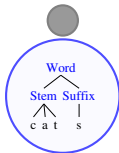
Suffix → Phoneme<sup>\*</sup>



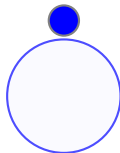
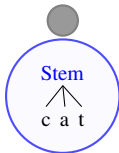
Generated words: cats

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

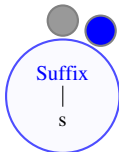
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



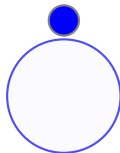
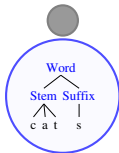
Suffix → Phoneme<sup>\*</sup>



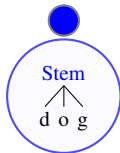
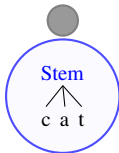
Generated words: cats

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

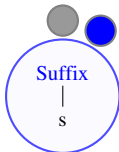
Word  $\rightarrow$  Stem Suffix



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



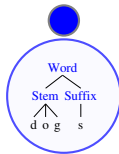
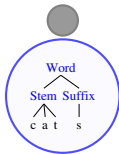
Suffix  $\rightarrow$  Phoneme<sup>\*</sup>



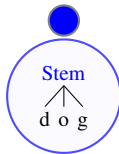
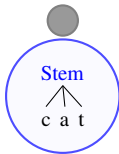
Generated words: cats

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

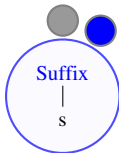
Word → Stem Suffix



Stem → Phoneme<sup>+</sup>



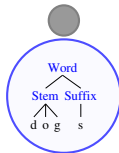
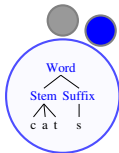
Suffix → Phoneme<sup>\*</sup>



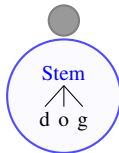
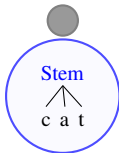
Generated words: cats, dogs

# Adaptor grammar for stem-suffix morphology (3)

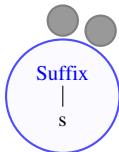
Word  $\rightarrow$  Stem Suffix



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



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



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

# Adaptor grammars as generative processes

- *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 nonterminal  $A$  has a CRP (or PYP) that caches previously generated subtrees rooted in  $A$
- The CFG rules of the adapted nonterminals determine the *base distributions* of these CRPs/PYPs
- The trees generated by an adaptor grammar are *not* independent
  - ▶ if an adapted subtree has been used frequently in the past, it's more likely to be used again

⇒ an adaptor grammar *learns* from the trees it generates
- but the sequence of trees is *exchangable* (important for sampling)

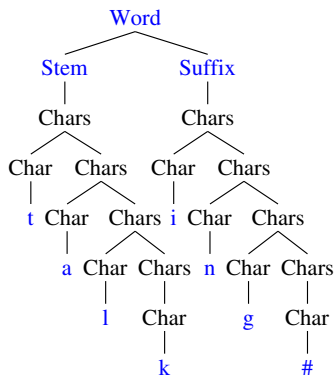


# Properties of adaptor grammars

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

# 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



# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

**Bayesian inference for adaptor grammars**

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# Estimating adaptor grammars

- Need to estimate:
  - ▶ parse trees  $\mathbf{t} = (t_1, \dots, t_n)$  for strings  $\mathbf{x} = (x_1, \dots, x_n)$
  - ▶ cached subtrees  $\tau$  for adapted nonterminals
  - ▶ DP parameters  $\alpha$  for adapted nonterminals
  - ▶ probabilities  $\theta$  of base grammar rules
- Collapsed Metropolis-within-Gibbs sampler for parse trees
  - ▶ sample parse tree  $t_i$  from  $P(t \mid x_i, \mathbf{t}_{-i})$   
where  $\mathbf{t}_{-i} = (t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n)$
  - ▶ sampling directly from conditional distribution of parses seems intractable (just as for PCFGs)
  - ▶ construct PCFG *proposal grammar*  $G'(\mathbf{t}_{-i})$  on the fly  
 $G'(\mathbf{t}_{-i})$  contains a rule for each cached subtree  $\tau$   
 $\Rightarrow$  grammar rules change while sampling
  - ▶ MH accept/reject determines whether to update  $t_i$

# Collapsed Metropolis-within-Gibbs sampling

- Input: terminal strings  $\mathbf{x} = (x_1, \dots, x_n)$
- Output: stream of sample parse trees  $\mathbf{t} = (t_1, \dots, t_n)$
- Metropolis-within-Gibbs sampling algorithm:
  - initialize  $\mathbf{t}$  somehow (e.g., random trees)
  - repeat forever:
    - pick a sentence  $x_i$  at random
    - construct PCFG proposal grammar  $G'(t_{-i})$  on the fly
    - sample parse tree  $t$  from  $P(t \mid x_i, G'(t_{-i}))$
    - update  $t_i$  with  $t$  according to MH accept-reject procedure
- After *burn-in* the samples  $\mathbf{t}$  are distributed according to  $P(\mathbf{t} \mid \mathbf{x})$

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

**Adaptor grammars for unsupervised word segmentation**

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# Unsupervised word segmentation

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

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

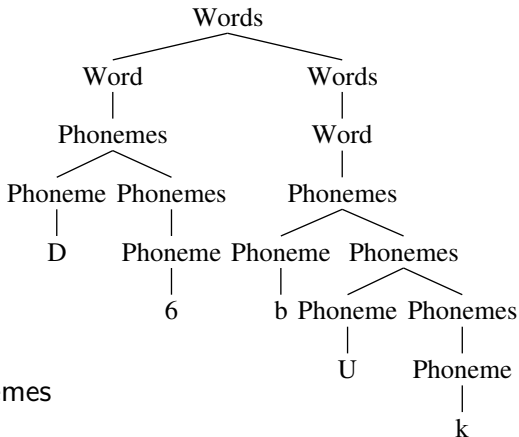
- Useful cues for word segmentation:
  - ▶ Phonotactics (Fleck)
  - ▶ Inter-word dependencies (Goldwater)

# 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  $\rightarrow$   $a \mid \dots \mid z$



- PCFG trees are adequate for representing word segmentation, but PCFG rules are *too small a unit of generalization* to learn words

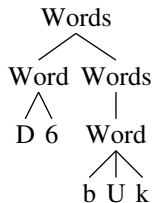


# Word segmentation with PCFGs (1)

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

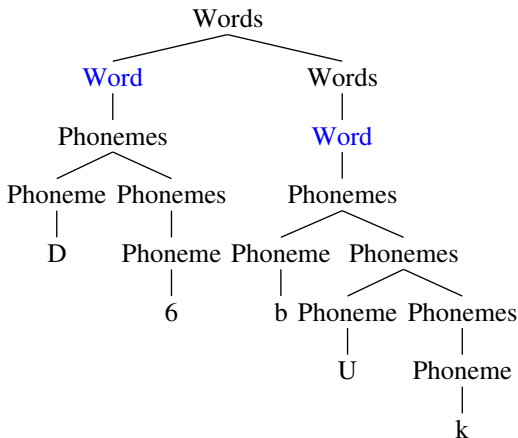


# Unigram word segmentation adaptor grammar

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*

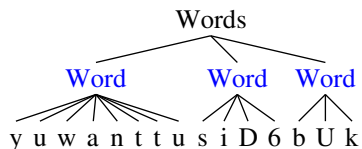
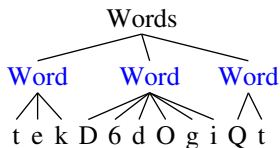
100 Word  $\rightarrow$  *D 6 d O g i*

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)

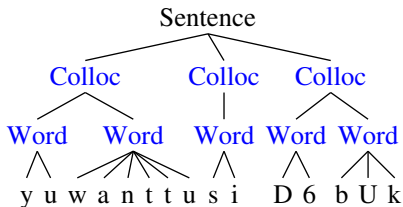


# colloc: Collocations $\Rightarrow$ Words

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>

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

Word  $\rightarrow$  SyllableIF

Syllable  $\rightarrow$  (Onset) Rhyme

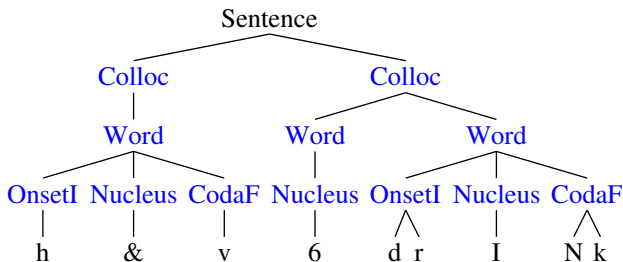
Word  $\rightarrow$  SyllableI (Syllable) (Syllable) SyllableF

Rhyme  $\rightarrow$  Nucleus (Coda)

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

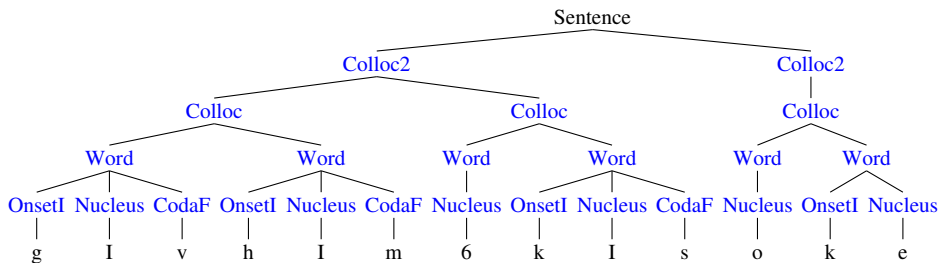
Coda  $\rightarrow$  Consonant<sup>+</sup>

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



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

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



# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

**Topic Collocation models using adaptor grammars**

Adaptor grammars for named entities

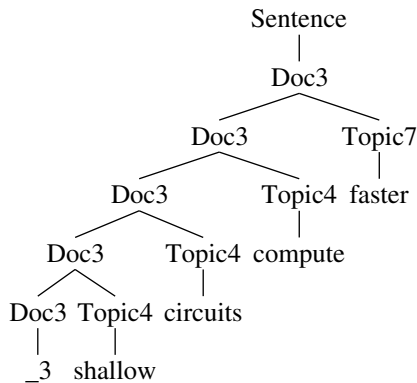
Conclusion

Extending Adaptor Grammars



# LDA topic models as PCFGs

- Each document  $i$  generates a distribution over  $m$  topics
- Each topic  $j$  generates a (unigram) distribution over vocabulary  $\mathcal{X}$ .
- Preprocess input by prepending a document id to every sentence



Sentence  $\rightarrow$  Doc $_i$        $i \in 1, \dots, n$   
Doc $_i \rightarrow$   $_i$        $i \in 1, \dots, n$   
Doc $_i \rightarrow$  Doc $_i$  Topic $_j$        $i \in 1, \dots, n; j \in 1, \dots, m$   
Topic $_j \rightarrow$   $x$        $j \in 1, \dots, m; x \in \mathcal{X}$

## Estimated PCFG for LDA topic model

Rule	Count
$\text{Doc}_3 \rightarrow \text{Doc}_3 \text{ Topic}_4$	737
$\text{Doc}_3 \rightarrow \text{Doc}_3 \text{ Topic}_3$	392
$\text{Doc}_3 \rightarrow \text{Doc}_3 \text{ Topic}_9$	263
$\text{Doc}_3 \rightarrow \text{Doc}_3 \text{ Topic}_7$	254
...	
$\text{Topic}_4 \rightarrow \text{function}$	3576
$\text{Topic}_4 \rightarrow \text{functions}$	2195
$\text{Topic}_4 \rightarrow \text{error}$	1614
$\text{Topic}_4 \rightarrow \text{algorithm}$	1588
...	
$\text{Topic}_7 \rightarrow \text{network}$	13018
$\text{Topic}_7 \rightarrow \text{input}$	6716
$\text{Topic}_7 \rightarrow \text{training}$	6078
$\text{Topic}_7 \rightarrow \text{learning}$	5176
...	

# Finding topic collocations with adaptor grammars

- Let each  $\text{Topic}_j$  expand to a *sequence* of words

- Adapt each  $\text{Topic}_j$

⇒ learn topic's likely collocations

Sentence  $\rightarrow$   $\text{Doc}_i$

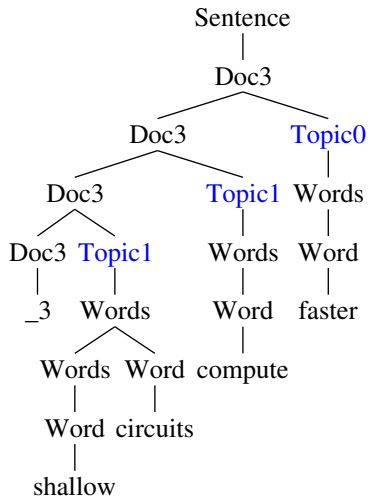
$\text{Doc}_i \rightarrow$   ${}_i$

$\text{Doc}_i \rightarrow \text{Doc}_i \text{Topic}_j$

Topic $_j \rightarrow$  Words

Words  $\rightarrow$  Word

Words  $\rightarrow$  Words Word



- Would be easy to make a hierarchical adaptor grammar model that shares collocations between topics

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

**Adaptor grammars for named entities**

Conclusion

Extending 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  $\rightarrow$  (A0) (A1) ... (A6)    NP  $\rightarrow$  (B0) (B1) ... (B6)

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

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

...

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

...

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

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

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

(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)

- See Elsner et al (2009) for details

# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

## Conclusion

Extending Adaptor Grammars

# Summary and future work

- Parsing and learning are intimately related
  - ▶ most methods for grammar learning involve repeated parsing
- *Non-parametric generalizations* of PCFGs introduce an infinite number of potential rules
  - ▶ *Adaptor Grammars* learn entire subtrees
  - ▶ can define a range of interesting models
- Collapsed Gibbs samplers are a natural method for non-parametric inference
  - ▶ a finite number of samples can only involve a finite number of structures
  - ▶ efficient sampling using Metropolis-within-Gibbs with PCFG proposals
- *Non-parametric Bayesian approaches provide a principled approach to learning grammar rules as well as their probabilities*



# Outline

Learning rule probabilities

Learning grammar rules (not just probabilities)

Chinese Restaurant Processes

Adaptor grammars

Bayesian inference for adaptor grammars

Adaptor grammars for unsupervised word segmentation

Topic Collocation models using adaptor grammars

Adaptor grammars for named entities

Conclusion

Extending Adaptor Grammars

# Issues with adaptor grammars

- Recursion *through adapted nonterminals* seems problematic
  - ▶ New tables are created as each node is encountered top-down
  - ▶ But the tree labeling the table is only known after the whole subtree has been completely generated
  - ▶ If adapted nonterminals are recursive, might pick a table whose label we are currently constructing. What then?
- Extend adaptor grammars so adapted fragments can end at nonterminals a la DOP (currently always go to terminals)
  - ▶ Adding “exit probabilities” to each adapted nonterminal
  - ▶ In some approaches, fragments can grow “above” existing fragments, but can’t grow “below” (O’Donnell)
- Adaptor grammars *conflate grammatical and Bayesian hierarchies*
  - ▶ Might be useful to disentangle them with *meta-grammars*

# Context-free grammars

A *context-free grammar* (CFG) consists of:

- a finite set  $N$  of *nonterminals*,
- a finite set  $W$  of *terminals* disjoint from  $N$ ,
- a finite set  $R$  of *rules*  $A \rightarrow \beta$ , where  $A \in N$  and  $\beta \in (N \cup W)^*$
- a *start symbol*  $S \in N$ .

Each  $A \in N \cup W$  *generates* a set  $\mathcal{T}_A$  of trees.

These are the smallest sets satisfying:

- If  $A \in W$  then  $\mathcal{T}_A = \{A\}$ .
- If  $A \in N$  then:

$$\mathcal{T}_A = \bigcup_{A \rightarrow B_1 \dots B_n \in R_A} \text{TREE}_A(\mathcal{T}_{B_1}, \dots, \mathcal{T}_{B_n})$$

where  $R_A = \{A \rightarrow \beta : A \rightarrow \beta \in R\}$ , and

$$\text{TREE}_A(\mathcal{T}_{B_1}, \dots, \mathcal{T}_{B_n}) = \left\{ \begin{array}{l} A \\ \wedge \\ t_1 \dots t_n \end{array} : \begin{array}{l} t_i \in \mathcal{T}_{B_i}, \\ i = 1, \dots, n \end{array} \right\}$$

The set of trees generated by a CFG is  $\mathcal{T}_S$ .

## Probabilistic context-free grammars

A *probabilistic context-free grammar* (PCFG) is a CFG and a vector  $\theta$ , where:

- $\theta_{A \rightarrow \beta}$  is the probability of expanding the nonterminal  $A$  using the production  $A \rightarrow \beta$ .

It defines distributions  $G_A$  over trees  $\mathcal{T}_A$  for  $A \in N \cup W$ :

$$G_A = \begin{cases} \delta_A & \text{if } A \in W \\ \sum_{A \rightarrow B_1 \dots B_n \in R_A} \theta_{A \rightarrow B_1 \dots B_n} \text{TD}_A(G_{B_1}, \dots, G_{B_n}) & \text{if } A \in N \end{cases}$$

where  $\delta_A$  puts all its mass onto the singleton tree  $A$ , and:

$$\text{TD}_A(G_1, \dots, G_n) \left( \begin{array}{c} A \\ \wedge \\ t_1 \dots t_n \end{array} \right) = \prod_{i=1}^n G_i(t_i).$$

$\text{TD}_A(G_1, \dots, G_n)$  is a distribution over  $\mathcal{T}_A$  where each subtree  $t_i$  is generated independently from  $G_i$ .

## DP adaptor grammars

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

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

$$H_A = \sum_{A \rightarrow B_1 \dots B_n \in R_A} \theta_{A \rightarrow B_1 \dots B_n} \text{TD}_A(G_{B_1}, \dots, G_{B_n})$$

The grammar generates the distribution  $G_S$ .

One Dirichlet Process for each adapted non-terminal  $A$  (i.e.,  $\alpha_A > 0$ ).

# Recursion in adaptor grammars

- The probability of joint distributions  $(G, H)$  is defined by:

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

$$H_A = \sum_{A \rightarrow B_1 \dots B_n \in R_A} \theta_{A \rightarrow B_1 \dots B_n} \text{TD}_A(G_{B_1}, \dots, G_{B_n})$$

- This holds *even if adaptor grammar is recursive*
- Question: when does this define a *distribution* over  $(G, H)$ ?

# Adaptive fragment grammars

- Disentangle syntactic and Bayesian hierarchy
  - ▶ *Adaptive metagrammar* generates *fragment distributions*
  - ▶ which plug together as in tree substitution grammar
- Tree fragment sets  $\mathcal{P}_A, A \in N$  are smallest sets satisfying:

$$\mathcal{P}_A = \bigcup_{A \rightarrow B_1 \dots B_n \in R_A} \text{TREE}_A(\{B_1\} \cup \mathcal{P}_{B_1}, \dots, \{B_n\} \cup \mathcal{P}_{B_n})$$

- Grammar's distributions  $G_A$  over  $\mathcal{T}_A$  defined using *fragment distributions*  $F_A$  over  $\mathcal{P}_A$  (*generalized PCFG rules*)

$$G_A = \sum_{\substack{\text{triangle} \\ B_1 \dots B_n \\ \in \mathcal{P}_A}} F_A(\text{triangle}) \text{TD}_A(\text{triangle}, G_{B_1}, \dots, G_{B_n})$$

- A fragment grammar generates the distribution  $G_S$

# Adaptive fragment distributions

- $H_A$  is a PCFG distribution over  $\mathcal{P}_A$

$$H_A = \sum_{A \rightarrow B_1 \dots B_n \in R_A} \theta_{A \rightarrow B_1 \dots B_n} \text{TD}_A(\eta \delta_{B_1} + (1 - \eta)H_{B_1}, \dots)$$

where  $\eta$  is the *fragment exit probability*

- Obtain  $F_A$  by *adapting the  $H_A$  distribution*

$$F_A \sim \text{DP}(\alpha_A, H_A)$$

- This construction can be *iterated*, i.e., replace  $\theta$  with another fragment distribution
- Question: if we iterate this, when does the *fixed point* exist, and what is it?