

Grammars and Topic Models

Mark Johnson
Macquarie University
Sydney
Australia

joint work with Michael Frank, Sharon Goldwater, Tom Griffiths and Bevan
Jones

December 2009

Talk summary

- Probabilistic Context-Free Grammars (PCFGs)
 - ▶ a way of specifying certain stochastic automata
- LDA topic models as PCFGs
- Non-parametric generalizations of PCFGs
 - ▶ “infinite” PCFGs
 - ▶ adaptor grammars
- Topic models and adaptor grammars
 - ▶ topical collocations
 - ▶ the structure of names
 - ▶ learning words and their meanings

Outline

Probabilistic Context-Free Grammars

LDA topic models as PCFGs

Adaptor grammars

Adaptor grammars and topic models

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 A with a multinomial distribution θ_A over the rules $R_A = \{A \rightarrow \alpha\}$ that expand it
- Probability of a tree is the *product of the probabilities of the rules* used to construct it

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

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

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

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

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 θ , where:

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

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 δ_A puts all its mass onto the singleton tree A , and:

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

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

Bayesian inference for rule probabilities

- A PCFG specifies a *multinomial* distribution θ_A for each nonterminal A over the rules $A \rightarrow \alpha$ expanding A
 - ▶ $\theta_{A \rightarrow \alpha}$ is probability of A expanding to α
 - ▶ probability of a parse t is a *product of multinomials*
- Conjugate prior: $\theta_A \mid \alpha_A \sim \text{Dir}(\alpha_A)$ for each nonterminal A
- Given a corpus of *parse trees* \mathbf{t} , posterior distribution for $\theta_A \mid \mathbf{t}, \alpha_A \sim \text{Dir}(\alpha_A + \mathbf{n}_A(\mathbf{t}))$
- Given a corpus of *strings* \mathbf{w} (the *yields* of \mathbf{t}), joint posterior distribution over
 - ▶ rule probabilities θ
 - ▶ parse trees \mathbf{t}

is intractable, but can be approximated by:

- ▶ variational Bayes (mean field approximation)
- ▶ Markov chain Monte Carlo

Outline

Probabilistic Context-Free Grammars

LDA topic models as PCFGs

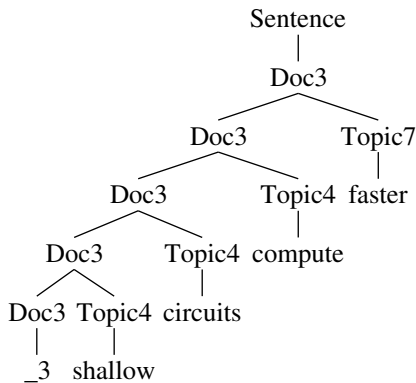
Adaptor grammars

Adaptor grammars and topic models

Conclusion

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 $_j \rightarrow x$ $j \in 1, \dots, m; x \in \mathcal{X}$

Left spine identifies document

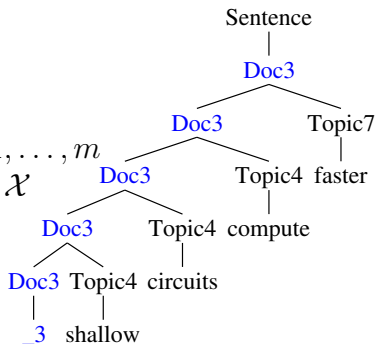
- Left spine passes document id throughout sentence

Sentence \rightarrow Doc_{*i*} *i* \in 1, ..., *n*

Doc_{*i*} \rightarrow _{*i*} *i* \in 1, ..., *n*

Doc_{*i*} \rightarrow Doc_{*i*} Topic_{*j*} *i* \in 1, ..., *n*; *j* \in 1, ..., *m*

Topic_{*j*} \rightarrow *x* *j* \in 1, ..., *m*; *x* \in \mathcal{X}



Document \rightarrow topic rules

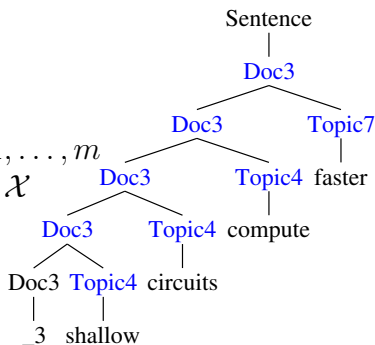
- Document \rightarrow topic rules specify probability of topic within document

Sentence \rightarrow Doc_{*i*} *i* \in 1, ..., *n*

Doc_{*i*} \rightarrow $_i$ *i* \in 1, ..., *n*

Doc_{*i*} \rightarrow Doc_{*i*} Topic_{*j*} *i* \in 1, ..., *n*; *j* \in 1, ..., *m*

Topic_{*j*} \rightarrow *x* *j* \in 1, ..., *m*; *x* \in \mathcal{X}



Topic \rightarrow word rules

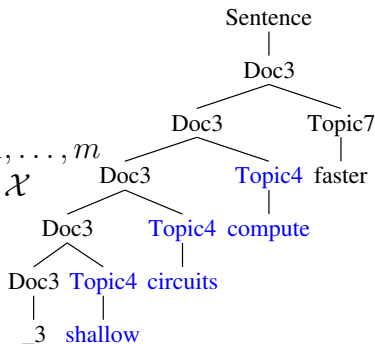
- Topic \rightarrow word rules specify probability of word within topic

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 $j \in 1, \dots, m; x \in \mathcal{X}$



Bayesian inference for LDA PCFGs

- Dirichlet priors on Document \rightarrow Topic and Topic \rightarrow Word distributions
- General-purpose PCFG parsing/estimation algorithms require time cubic in length of sentence
 - ▶ not a good idea for long documents!
- More efficient algorithms for these kinds of grammars
 - ▶ (standard LDA inference algorithms)
 - ▶ predictive (e.g., Earley) parsing algorithms
 - ▶ identify compositions finite state automata/transducers that these grammars encode

Extended LDA PCFG (1): sticky topics

- HMM generating *sequences of topics* in each document
- Non-uniform Dirichlet prior over topic \rightarrow topic transitions \Rightarrow sticky topics
- Grammar rule schemata given
 - ▶ document identifiers \mathcal{D} ,
 - ▶ topics $\mathcal{T} = \{1, \dots, m\}$, and
 - ▶ vocabulary \mathcal{W}

Sentence \rightarrow Doc $_{d,t}$

for each $d \in \mathcal{D}$ and $t \in \mathcal{T}$

Doc $_{d,t} \rightarrow d$ Topic $_t$

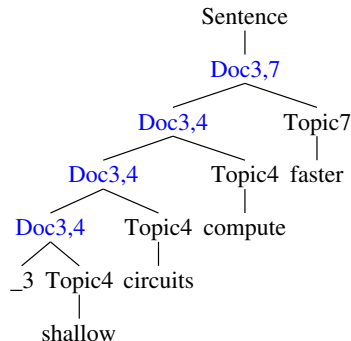
for each $d \in \mathcal{D}$ and $t \in \mathcal{T}$

Doc $_{d,t} \rightarrow$ Doc $_{d,t'}$ Topic $_t$

for each $d \in \mathcal{D}$ and $t, t' \in \mathcal{T}$

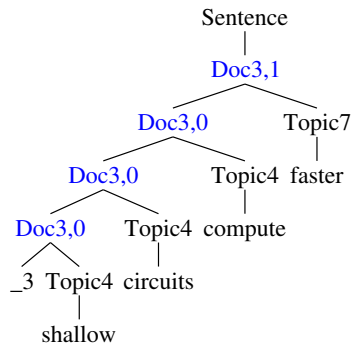
Topic $_t \rightarrow w$

for each $t \in \mathcal{T}$ and $w \in \mathcal{W}$



Extended LDA PCFG (2): document segmentation

- Divide documents into *segments*
- Grammar rule schemata given
 - ▶ document identifiers \mathcal{D} ,
 - ▶ *segments* $\mathcal{S} = \{1, \dots, \ell\}$,
 - ▶ topics $\mathcal{T} = \{1, \dots, m\}$, and
 - ▶ vocabulary \mathcal{W}



Sentence \rightarrow Doc $_{d,s}$

for each $d \in \mathcal{D}$ and $s \in \mathcal{S}$

Doc $_{d,s} \rightarrow d$

for each $d \in \mathcal{D}$ and $s \in \mathcal{S}$

Doc $_{d,s} \rightarrow$ Doc $_{d,s}$ Topic $_t$

for each $d \in \mathcal{D}$, $s \in \mathcal{S}$ and $t \in \mathcal{T}$

Doc $_{d,s} \rightarrow$ Doc $_{d,s'}$

for each $d \in \mathcal{D}$ and $s, s' \in \mathcal{S}$

Topic $_t \rightarrow w$

for each $t \in \mathcal{T}$ and $w \in \mathcal{W}$

What's the point of using grammars?

- All of these models can be stated directly (without using grammars)
 - ▶ grammars don't make inference easier or produce better results
- Grammars provide another way of *describing complex models*
 - ▶ mutually recursive hierarchies of sequential structures
- There are generic algorithms for *PCFG parsing and inference*
 - ▶ rapid proto-typing of new models
 - ▶ may lead to (more) efficient implementation via *stochastic automata*

Outline

Probabilistic Context-Free Grammars

LDA topic models as PCFGs

Adaptor grammars

Adaptor grammars and topic models

Conclusion

Nonparametric generalizations of PCFGs

- Two obvious nonparametric extensions of PCFGs:
 - ▶ let the number of nonterminals N 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 R 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 ...

A CFG for stem-suffix morphology

Word \rightarrow Stem Suffix

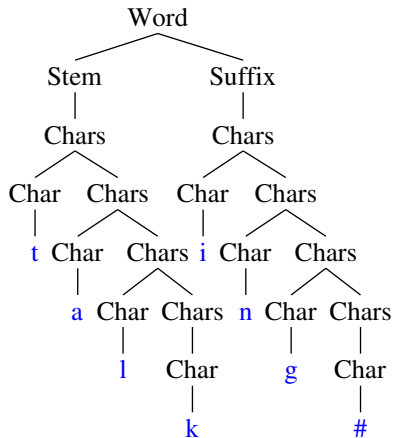
Stem \rightarrow Chars

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

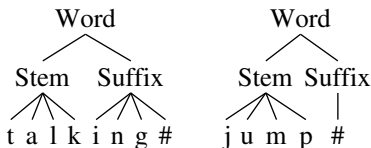
- But grammar's *units of generalization (PCFG rules)* are "too small" to learn morphemes

A “CFG” with one rule per possible morpheme

Word \rightarrow Stem Suffix

Stem \rightarrow *all possible stems*

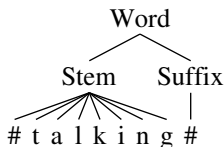
Suffix \rightarrow *all possible suffixes*



- 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

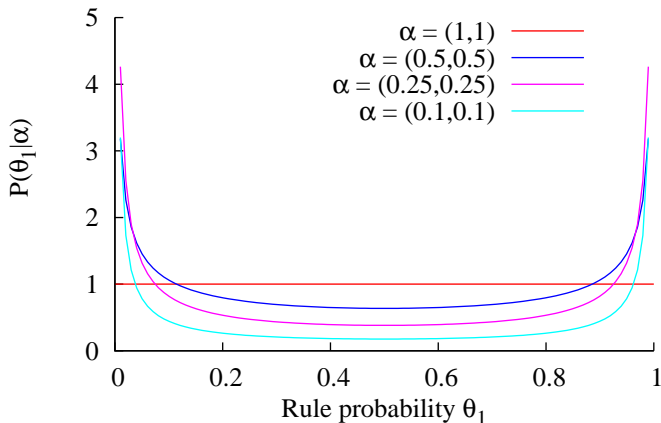
Maximum likelihood estimate for θ is trivial

- Maximum likelihood selects θ 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 in standard PCFG estimation, but can “turn rule off” by setting $\theta_{A \rightarrow \beta} \approx 0$
- Dirichlet prior with $\alpha_{A \rightarrow \beta} \approx 0$ prefers $\theta_{A \rightarrow \beta} \approx 0$



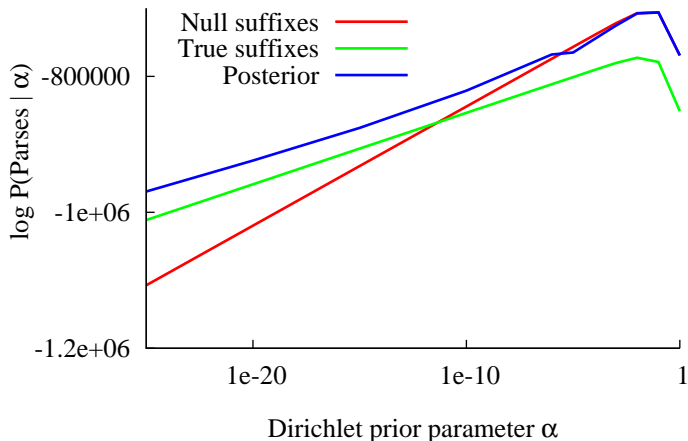
Morphological segmentation experiment

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

Posterior samples from WSJ verb tokens

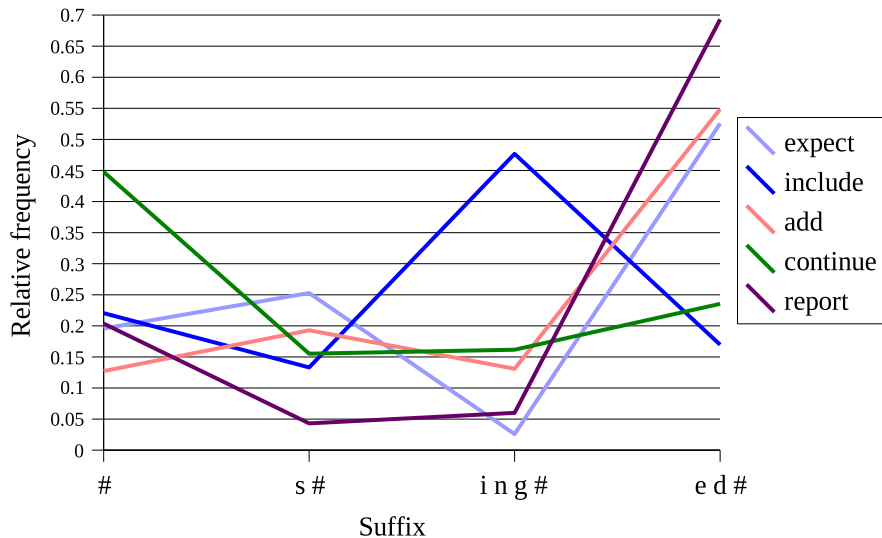
$\alpha = 0.1$	$\alpha = 10^{-5}$	$\alpha = 10^{-10}$	$\alpha = 10^{-15}$
expect	expect	expect	expect
expects	expects	expects	expects
expected	expected	expected	expected
expecting	expect ing	expect ing	expect ing
include	include	include	include
includes	includes	includ es	includ es
included	included	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 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”

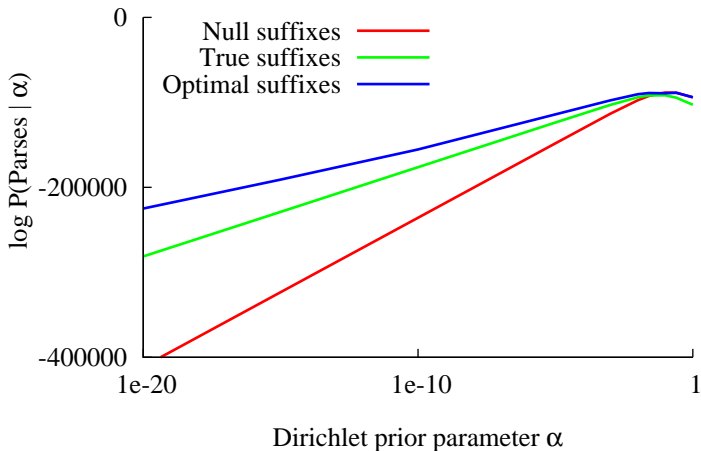
Types = “the”, “cat”, “chased”, “other”

- Estimating θ 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
- Adaptor grammar mimics Goldwater et al “Interpolating between Types and Tokens” morphology-learning model

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

Adaptor grammars: informal description

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

Adaptor grammar for stem-suffix morphology (0)

Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



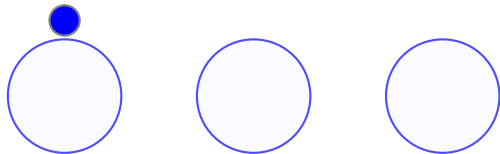
Suffix \rightarrow Phoneme^{*}



Generated words:

Adaptor grammar for stem-suffix morphology (1a)

Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



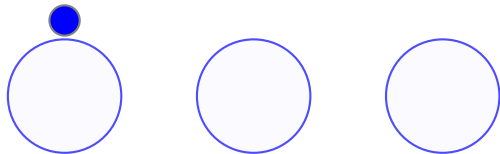
Suffix \rightarrow Phoneme^{*}



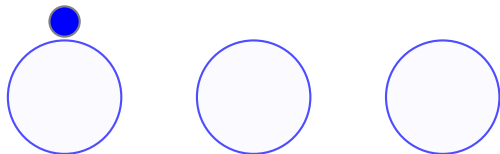
Generated words:

Adaptor grammar for stem-suffix morphology (1b)

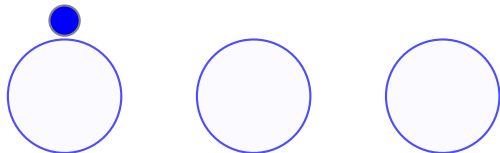
Word → Stem Suffix



Stem → Phoneme⁺



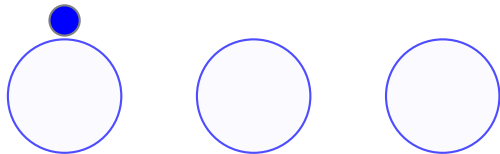
Suffix → Phoneme^{*}



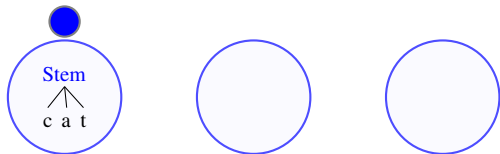
Generated words:

Adaptor grammar for stem-suffix morphology (1c)

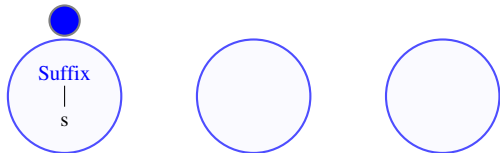
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



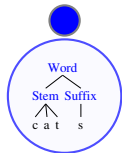
Suffix \rightarrow Phoneme^{*}



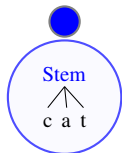
Generated words:

Adaptor grammar for stem-suffix morphology (1d)

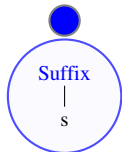
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



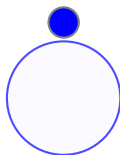
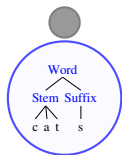
Suffix \rightarrow Phoneme^{*}



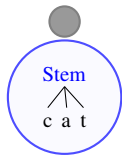
Generated words: *cats*

Adaptor grammar for stem-suffix morphology (2a)

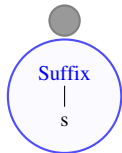
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



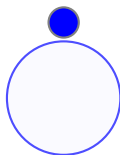
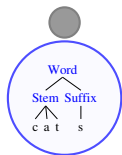
Suffix \rightarrow Phoneme^{*}



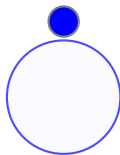
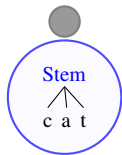
Generated words: cats

Adaptor grammar for stem-suffix morphology (2b)

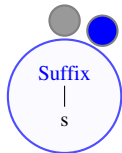
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



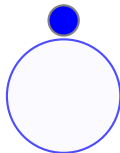
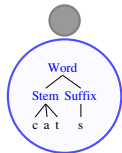
Suffix \rightarrow Phoneme^{*}



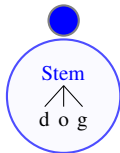
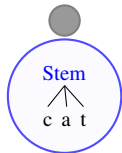
Generated words: cats

Adaptor grammar for stem-suffix morphology (2c)

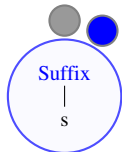
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



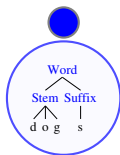
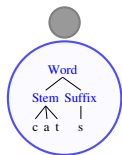
Suffix \rightarrow Phoneme^{*}



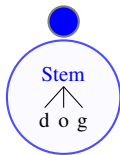
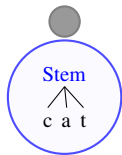
Generated words: cats

Adaptor grammar for stem-suffix morphology (2d)

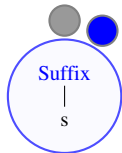
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



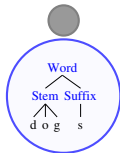
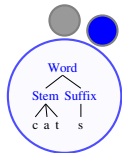
Suffix \rightarrow Phoneme^{*}



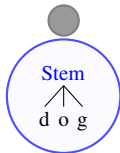
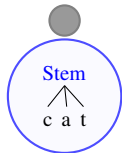
Generated words: cats, dogs

Adaptor grammar for stem-suffix morphology (3)

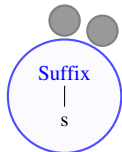
Word \rightarrow Stem Suffix



Stem \rightarrow Phoneme⁺



Suffix \rightarrow Phoneme^{*}



Generated words: cats, dogs, *cats*

Adaptor grammars as generative processes

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

DP adaptor grammars

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

$$\begin{aligned} G_A &\sim \text{DP}(\alpha_A, H_A) \text{ if } \alpha_A > 0, \text{ i.e., } A \text{ is } \textit{adapted} \\ &= H_A \quad \text{otherwise} \end{aligned}$$

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

The grammar generates the distribution G_S .

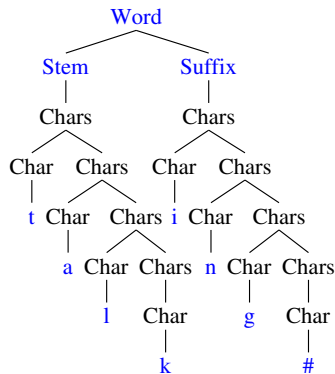
One Dirichlet Process for each adapted non-terminal A .

Properties of adaptor grammars

- Possible trees are generated by CFG rules
but the probability of each adapted tree is learned separately
 - Probability of adapted subtree τ is proportional to:
 - ▶ the number of times τ was seen before
⇒ “rich get richer” dynamics (Zipf distributions)
 - ▶ plus α_A times prob. of generating it via PCFG expansion
- ⇒ Useful compound structures can be *more probable than their parts*
- PCFG rule probabilities estimated *from table labels*
 - ⇒ effectively *learns from types*, not tokens
 - ⇒ makes learner less sensitive to frequency variation in input

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



Unsupervised word segmentation

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

y_Δu_Δw_Δa_Δn_Δt_Δt_Δu_Δs_Δi_ΔD_Δ6_Δb_ΔU_Δk

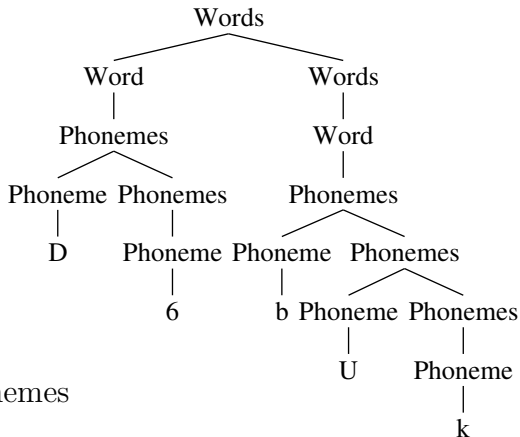
- Useful cues for word segmentation:
 - ▶ Phonotactics (Fleck)
 - ▶ Inter-word dependencies (Goldwater)

Word segmentation with PCFGs (1)

Sentence \rightarrow Word⁺
Word \rightarrow Phoneme⁺

which abbreviates

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

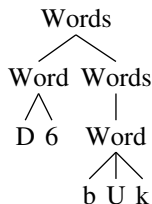


Word segmentation with PCFGs (1)

Sentence \rightarrow Word⁺

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

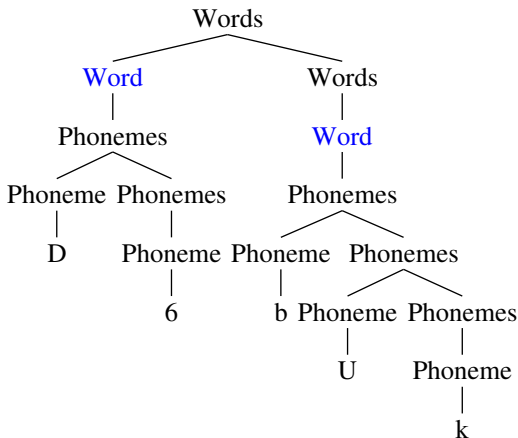


Unigram word segmentation adaptor grammar

Sentence \rightarrow Word⁺

Word \rightarrow Phoneme⁺

- *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

Adaptor grammar learnt from Brent corpus

- Initial grammar*

1	Sentence \rightarrow Word Sentence	1	Sentence \rightarrow Word
1	Word \rightarrow Phons		
1	Phons \rightarrow Phon Phons	1	Phons \rightarrow Phon
1	Phon \rightarrow D	1	Phon \rightarrow G
1	Phon \rightarrow A	1	Phon \rightarrow E

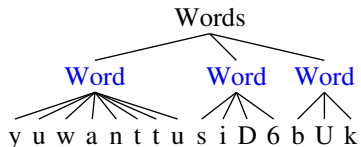
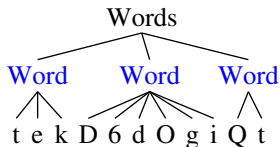
- A grammar learnt from Brent corpus*

16625	Sentence \rightarrow Word Sentence	9791	Sentence \rightarrow Word
1	Word \rightarrow Phons		
4962	Phons \rightarrow Phon Phons	1575	Phons \rightarrow Phon
134	Phon \rightarrow D	41	Phon \rightarrow G
180	Phon \rightarrow A	152	Phon \rightarrow E
460	Word \rightarrow (Phons (Phon <i>y</i>) (Phons (Phon <i>u</i>)))		
446	Word \rightarrow (Phons (Phon <i>w</i>) (Phons (Phon <i>A</i>) (Phons (Phon <i>t</i>))))		
374	Word \rightarrow (Phons (Phon <i>D</i>) (Phons (Phon <i>6</i>)))		
372	Word \rightarrow (Phons (Phon <i>&</i>) (Phons (Phon <i>n</i>) (Phons (Phon <i>d</i>))))		

Words (unigram model)

Sentence \rightarrow Word⁺ Word \rightarrow Phoneme⁺

- 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)

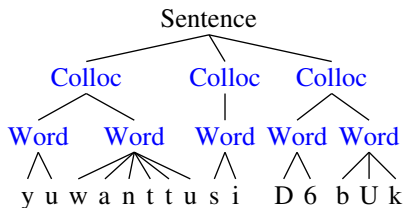


Collocations \Rightarrow Words

Sentence \rightarrow Colloc⁺

Colloc \rightarrow Word⁺

Word \rightarrow Phon⁺



- 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)

Collocations \Rightarrow Words \Rightarrow Syllables

Sentence \rightarrow Colloc⁺

Word \rightarrow Syllable

Word \rightarrow Syllable Syllable Syllable

Onset \rightarrow Consonant⁺

Nucleus \rightarrow Vowel⁺

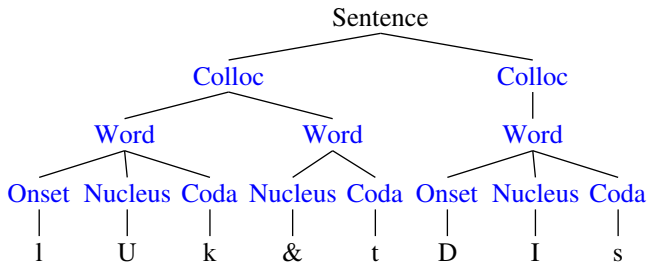
Colloc \rightarrow Word⁺

Word \rightarrow Syllable Syllable

Syllable \rightarrow (Onset) Rhyme

Rhyme \rightarrow Nucleus (Coda)

Coda \rightarrow Consonant⁺



- With no supra-word generalizations, f-score = 68%
- With 2 Collocation levels, f-score = 82%

Distinguishing internal onsets/codas helps

Sentence \rightarrow Colloc⁺

Word \rightarrow SyllableIF

Word \rightarrow SyllableI Syllable SyllableF

OnsetI \rightarrow Consonant⁺

Nucleus \rightarrow Vowel⁺

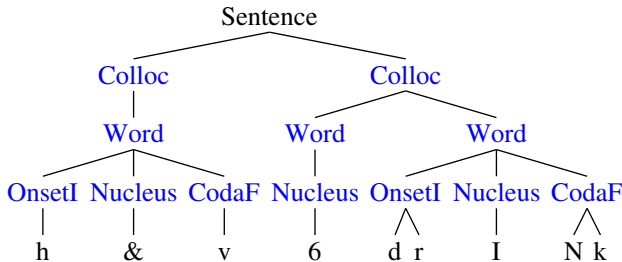
Colloc \rightarrow Word⁺

Word \rightarrow SyllableI SyllableF

SyllableIF \rightarrow (OnsetI) RhymeF

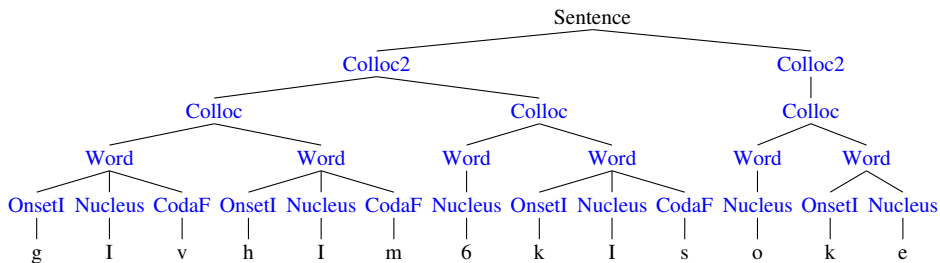
RhymeF \rightarrow Nucleus (CodaF)

CodaF \rightarrow Consonant⁺



- Without distinguishing initial/final clusters, f-score = 82%
- Distinguishing initial/final clusters, f-score = 84%
- With 2 Collocation levels, f-score = 87%

Collocations² ⇒ Words ⇒ Syllables

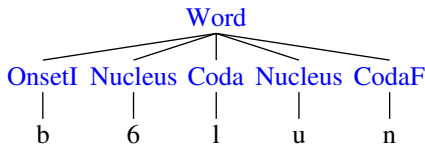


Syllabification learnt by adaptor grammars

- Grammar has no reason to prefer to parse word-internal intervocalic consonants as onsets

1 Syllable \rightarrow Onset Rhyme 1 Syllable \rightarrow Rhyme

- The learned grammars consistently analyse them as either Onsets or Codas \Rightarrow learns wrong grammar half the time

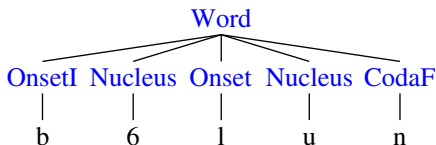


- Syllabification accuracy is relatively poor
Syllabification given true word boundaries: f-score = 83%
Syllabification learning word boundaries: f-score = 74%

Preferring Onsets improves syllabification

2 Syllable \rightarrow Onset Rhyme 1 Syllable \rightarrow Rhyme

- Changing the prior to prefer word-internal Syllables with Onsets dramatically improves segmentation accuracy
- “Rich get richer” property of Chinese Restaurant Processes
 \Rightarrow all ambiguous word-internal consonants analysed as Onsets



- Syllabification accuracy is much higher than without bias
Syllabification given true word boundaries: f-score = 97%
Syllabification learning word boundaries: f-score = 90%

Modelling sonority classes improves syllabification

$\text{Onset} \rightarrow \text{Onset}_{\text{Stop}}$	$\text{Onset} \rightarrow \text{Onset}_{\text{Fricative}}$
$\text{Onset}_{\text{Stop}} \rightarrow \text{Stop}$	$\text{Onset}_{\text{Stop}} \rightarrow \text{Stop Onset}_{\text{Fricative}}$
$\text{Stop} \rightarrow \text{p}$	$\text{Stop} \rightarrow \text{t}$

- Five consonant sonority classes
- $\text{Onset}_{\text{Stop}}$ generates a consonant cluster with a Stop at left edge
- Prior prefers transitions compatible with sonority hierarchy (e.g., $\text{Onset}_{\text{Stop}} \rightarrow \text{Stop Onset}_{\text{Fricative}}$) to transitions that aren't (e.g., $\text{Onset}_{\text{Fricative}} \rightarrow \text{Fricative Onset}_{\text{Stop}}$)
- Same transitional probabilities used for initial and non-initial Onsets (maybe not a good idea for English?)
- Word-internal Onset bias still necessary
- Syllabification given true boundaries: f-score = 97.5%
- Syllabification learning word boundaries: f-score = 91%

Summary: Adaptor grammars for word segmentation

- Easy to define adaptor grammars that are sensitive to:

Generalization	Accuracy
words as units (unigram)	56%
+ associations between words (collocations)	76%
+ syllable structure	87%

- word segmentation *improves when you learn other things as well*
 - ▶ *explain away* potentially misleading generalizations

Outline

Probabilistic Context-Free Grammars

LDA topic models as PCFGs

Adaptor grammars

Adaptor grammars and topic models

Conclusion

Finding topical collocations

- Modify LDA PCFG so each topic generates *sequences of words*
 - ▶ c.f. unigram word segmentation model

Sentence \rightarrow Doc _{d} for each $d \in \mathcal{D}$

Doc _{d} \rightarrow d for each $d \in \mathcal{D}$

Doc _{d} \rightarrow Doc _{d} Topic _{t} for each $d \in \mathcal{D}$ and $t \in \mathcal{T}$

Topic _{t} \rightarrow w^+ for each $t \in \mathcal{T}$ and $w \in \mathcal{W}$

where:

- ▶ document identifiers \mathcal{D} ,
- ▶ *topics* $\mathcal{T} = \{1, \dots, m\}$, and
- ▶ *vocabulary* \mathcal{W}

Sample output on NIPS corpus, 20 topics

- Multiword subtrees learned by adaptor grammar:

T_0 → gradient descent	T_1 → associative memory
T_0 → cost function	T_1 → standard deviation
T_0 → fixed point	T_1 → randomly chosen
T_0 → learning rates	T_1 → hamming distance
T_3 → membrane potential	T_10 → ocular dominance
T_3 → action potentials	T_10 → visual field
T_3 → visual system	T_10 → nervous system
T_3 → primary visual cortex	T_10 → action potential
- Sample parses:
 - _3 (T_5 polynomial size) (T_15 threshold circuits)
 - _4 (T_11 studied) (T_19 pattern recognition algorithms)
 - _4 (T_2 feedforward neural network) (T_1 implements)
 - _5 (T_11 single) (T_10 ocular dominance stripe) (T_12 low) (T_3

Learning the structure of 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.
 - ▶ used as part of a *coreference system*
- 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

Adaptor grammar for names

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

A0 \rightarrow Word⁺

B0 \rightarrow Word⁺

...

A6 \rightarrow Word⁺

...

B6 \rightarrow Word⁺

NP \rightarrow Unordered⁺

Unordered \rightarrow Word⁺

- *Sample parses:*

(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)

Learning words and their referents



PIG|DOG I_Δz_ΔD_Δ&_Δt_ΔD_Δ6_Δp_ΔI_Δg_Δ
PIG

- Input: unsegmented phoneme sequence and *objects in nonlinguistic context*
- Goal: segment into words and *learn word-object relationship*

AG for (unigram) segmentation and reference

- Given *possible referents* \mathcal{R} , the grammar contains rules:

Sentence	\rightarrow Referent _r	for each $r \in \mathcal{R}$
Referent _r	\rightarrow s	for each $r \in \mathcal{R}, r \in s \in 2^{\mathcal{R}}$
Referent _r	\rightarrow Referent _r Word _r	for each $r \in \mathcal{R}$
Referent _r	\rightarrow Referent _r Word _{\emptyset}	for each $r \in \mathcal{R}$
<u>Word</u> _r	\rightarrow Phonemes _r	for each $r \in \mathcal{R} \cup \{\emptyset\}$
Phonemes _r	\rightarrow Phoneme ⁺	for each $r \in \mathcal{R} \cup \{\emptyset\}$

- Sample parses:

T_dog|pig (Word I z D & t) (Word_dog D 6 d O g i)

T_dog|pig (Word D E r z) (Word_dog D 6 d O g i)

T_dog|pig (Word D & t s 6) (Word_pig p I g)

Joint segmentation and reference results

- Simultaneously learning word segmentation and reference does not seem to improve word segmentation
 - ▶ non-linguistic context is very impoverished
 - ▶ relatively few words are referential
- Referential words are segmented better when referents are provided
- Referential words are segmented better when at utterance ends
 - ▶ consistent with Frank's artificial language learning experiments

Outline

Probabilistic Context-Free Grammars

LDA topic models as PCFGs

Adaptor grammars

Adaptor grammars and topic models

Conclusion

Conclusion

- Grammars provide an alternative way of formulating complex models
- General-purpose inference procedures
- LDA topic models are PCFGs with Dirichlet priors on rule probabilities
- Two non-parametric generalizations of PCFGs
 - ▶ split non-terminals (states)
 - ▶ extend the set of possible rules (adaptor grammars)
- Adaptor grammars and topic models
 - ▶ topical collocations
 - ▶ structure in named entities
 - ▶ learning referents of words