

Unsupervised phonemic Chinese word segmentation using Adaptor Grammars

Mark Johnson¹ and Katherine Demuth²

¹Department of Computing

²Department of Linguistics

Macquarie University
Sydney
Australia

COLING, August 2010

Talk outline

- Adaptor grammars are a framework for expressing *non-parametric hierarchical Bayesian models*
- They can be used to define *unsupervised word segmentation* models that learn:
 - ▶ *word-internal structure*: how words are composed out of syllables, and
 - ▶ *inter-word structure*: collocational dependencies between words
- Adaptor Grammars provide state-of-the-art unsupervised segmentation results for English: *will they work for Mandarin Chinese?*
 - ▶ can Adaptor Grammars model *lexical tone*?
 - ▶ does modelling lexical tone *improve word segmentation accuracy*?

Why study computational models of language acquisition?

- Hypothesis: acquisition, comprehension and production are *computational processes*
 - ▶ computational models need not be just *descriptions* of language acquisition
 - ▶ a computational model should be able to *learn a language*
- Characterising computational models of acquisition:
 - ▶ the input (information available to learner)
 - ▶ the output (generalisations learner can make)
 - ▶ the algorithm used to map input to output
- *Bayesian inference algorithms* are optimal learners
 - ▶ computational generalisation of “ideal observer” theory
- Computational models let us study the effect of
 - ▶ changing the information in the input, and
 - ▶ altering the kinds of generalisations the learner can acquire in ways that would be impractical or unethical with real children
- *May be useful for designing experiments or therapeutic interventions*



Unsupervised word segmentation

- Input: phoneme sequences with *sentence boundaries* (Brent)
 - ▶ English data produced from orthographic transcripts of child-directed speech by *looking up each word in a pronouncing dictionary*
- Task: identify *word boundaries*, and hence words, in unsegmented utterance (in ARPABET)

y_Δu_▲w_Δa_Δn_Δt_▲t_Δu_▲s_Δi_▲D_Δ6_▲b_ΔU_Δk

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

CFG models of word segmentation

Words \rightarrow Word

Words \rightarrow Word Words

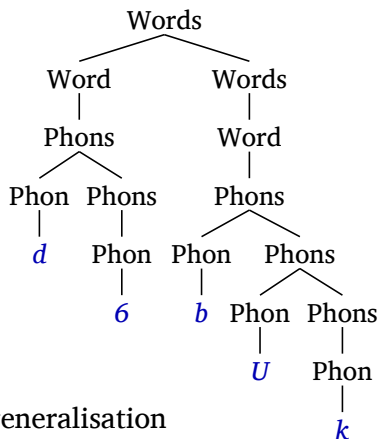
Word \rightarrow Phons

Phons \rightarrow Phon

Phons \rightarrow Phon Phons

Phon $\rightarrow a | b | \dots$

- CFG trees can *describe* segmentation, but
- PCFGs *can't distinguish* good segmentations from bad ones
 - ▶ PCFG rules are *too small* a unit of generalisation
 - ▶ need to learn e.g., probability that *bUk* is a Word



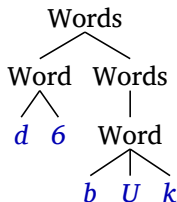
Towards non-parametric grammars

Words \rightarrow Word

Words \rightarrow Word Words

Word \rightarrow *all possible phoneme sequences*

- Learn probability Word \rightarrow b U k
- But *infinitely many possible Word expansions*
 \Rightarrow this grammar is *not a PCFG*
- Given *fixed training data*, only finitely many useful rules
 \Rightarrow use data to choose Word rules as well as their probabilities
- Non-parametric models: parameters of model depend on data



From PCFGs to Adaptor grammars

- An adaptor grammar is a PCFG where a subset of the nonterminals are *adapted*
- **Adaptor grammar generative process:**
 - ▶ to expand an *unadapted nonterminal* B : (just as in PCFG)
 - select a *rule* $B \rightarrow \beta \in R$ with prob. $\theta_{B \rightarrow \beta}$, and recursively expand nonterminals in β
 - ▶ to expand an *adapted nonterminal* B :
 - select a *previously generated subtree* T_B with prob. \propto number of times T_B was generated, or
 - select a *rule* $B \rightarrow \beta \in R$ with prob. $\propto \alpha_B \theta_{B \rightarrow \beta}$, and recursively expand nonterminals in β

Unigram adaptor grammar (Brent)

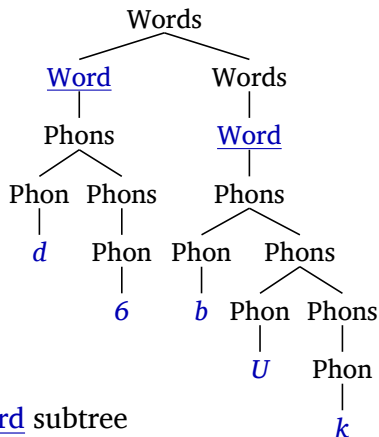
Words \rightarrow Word

Words \rightarrow Word Words

Word \rightarrow Phons

Phons \rightarrow Phon

Phons \rightarrow Phon Phons



- Word nonterminal is adapted

\Rightarrow To generate a Word:

- ▶ select a previously generated Word subtree with prob. \propto number of times it has been generated
- ▶ expand using Word \rightarrow Phons rule with prob. $\propto \alpha_{\text{Word}}$ and recursively expand Phons

Properties of adaptor grammars

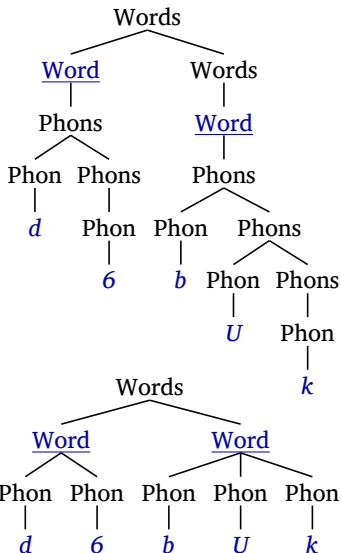
- Probability of regenerating an adapted subtree T_B
 \propto number of times T_B was previously generated
 - ▶ adapted subtrees are *not independent*
 - an adapted subtree can be *more probable* than the rules used to construct it
 - ▶ but they are *exchangable* \Rightarrow efficient sampling algorithms
 - ▶ “rich get richer” \Rightarrow Zipf power-law distributions
- Each adapted nonterminal is associated with a *Chinese Restaurant Process* or *Pitman-Yor Process*
 - ▶ CFG rules define *base distribution* of CRP or PYP
- CRP/PYP parameters (e.g., α_B) can themselves be estimated (e.g., slice sampling)

Abbreviatory notation

Words \rightarrow Word
Words \rightarrow Word Words
Word \rightarrow Phons
Phons \rightarrow Phon
Phons \rightarrow Phon Phons

is abbreviated as

Words \rightarrow Word⁺
Word \rightarrow Phon⁺

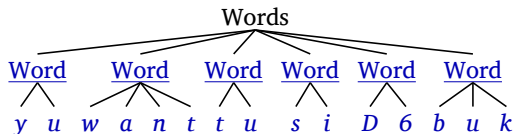


Unigram model of word segmentation

- Unigram “bag of words” model (Brent):
 - ▶ generate a *dictionary*, i.e., a set of words, where each word is a random sequence of phonemes
 - Bayesian prior prefers smaller dictionaries
 - ▶ generate each utterance by choosing each word at random from dictionary
- Brent’s unigram model as an Adaptor Grammar

Words \rightarrow Word⁺

Word \rightarrow Phon⁺



- Accuracy of word segmentation learnt: *56% token f-score* (same as Brent model)
- But we can construct many more word segmentation models using AGs

Adaptor grammar learnt from Brent corpus

- Initial grammar

1	Words \rightarrow <u>Word</u> Words	1	Words \rightarrow <u>Word</u>
1	<u>Word</u> \rightarrow Phon		
1	Phons \rightarrow Phon Phons	1	Phons \rightarrow Phon
1	Phon $\rightarrow D$	1	Phon $\rightarrow G$
1	Phon $\rightarrow A$	1	Phon $\rightarrow E$

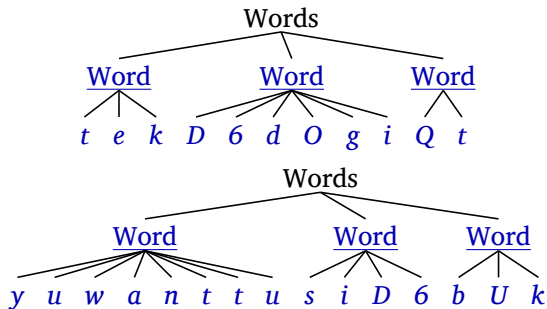
- A grammar learnt from Brent corpus

16625	Words \rightarrow <u>Word</u> Words	9791	Words \rightarrow <u>Word</u>
1575	<u>Word</u> \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	<u>Word</u> \rightarrow (Phons (Phon <i>y</i>) (Phons (Phon <i>u</i>)))		
446	<u>Word</u> \rightarrow (Phons (Phon <i>w</i>) (Phons (Phon <i>A</i>) (Phons (Phon <i>t</i>))))		
374	<u>Word</u> \rightarrow (Phons (Phon <i>D</i>) (Phons (Phon <i>6</i>)))		
372	<u>Word</u> \rightarrow (Phons (Phon <i>&</i>) (Phons (Phon <i>n</i>) (Phons (Phon <i>d</i>))))		

Undersegmentation errors with Unigram model

Words \rightarrow Word⁺ Word \rightarrow Phon⁺

- 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)

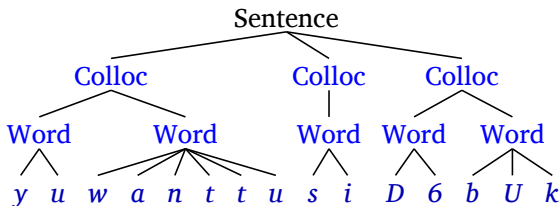


Collocations \Rightarrow Words

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 (74% 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⁺

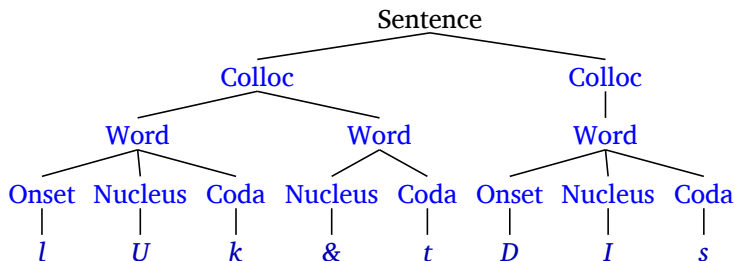
Colloc \rightarrow Word⁺

Word \rightarrow Syllable Syllable

Syllable \rightarrow (Onset) Rhyme

Rhyme \rightarrow Nucleus (Coda)

Coda \rightarrow Consonant⁺



- Rudimentary syllable model (an improved model might do better)
- With 2 Collocation levels, f-score = 84%

Distinguishing internal onsets/codas helps in English

Sentence \rightarrow Colloc⁺

Word \rightarrow SyllableIF

Word \rightarrow SyllableI Syllable SyllableF

OnsetI \rightarrow Consonant⁺

Nucleus \rightarrow Vowel⁺

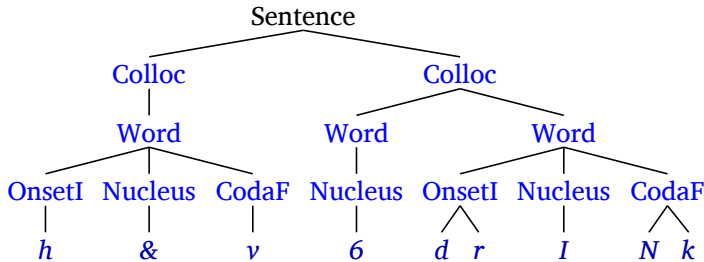
Colloc \rightarrow Word⁺

Word \rightarrow SyllableI SyllableF

SyllableIF \rightarrow (OnsetI) RhymeF

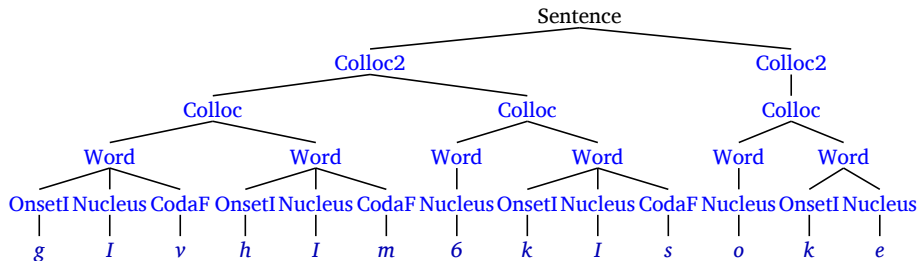
RhymeF \rightarrow Nucleus (CodaF)

CodaF \rightarrow Consonant⁺



- With 2 Collocation levels, not distinguishing initial/final clusters, f-score = 84%
- With 3 Collocation levels, distinguishing initial/final clusters,

Collocations² ⇒ Words ⇒ Syllables



Summary so far

- Word segmentation accuracy depends on the kinds of generalisations learnt.

Generalization	Accuracy
words as units (unigram)	56%
+ associations between words (collocations)	79%
+ syllable structure	87%

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

Tone in Mandarin Chinese word segmentation

- Tone in Mandarin Chinese provides an additional dimension of information to the language learner
- It is necessary in order to distinguish lexical items, but how important is it for word segmentation?
- Approach:
 - ▶ construct a pair of otherwise identical corpora, one that contains tone and one that does not
 - ▶ run identical learning algorithms on both corpora
 - ▶ compare the accuracy with which each learns word segmentation

Mandarin Chinese corpus

- Used Tardif (1993) “Beijing” corpus (in Pinyin format)
 - ▶ deleted all “Child” utterances, and utterances with codes \$INTERJ, \$UNINT, \$VOC and \$PRMPT
 - ▶ corpus contains 50,118 utterances, consisting of 187,533 word tokens

zen3me gei3 ta1 bei1 shang4 lai2 (1.) ?

ta1: (.) a1yi2 gei3 de (.) ta1 gei3 de .

hen3 jian3dan1 .

- Used a Pinyin to IPA translation program to produce IPA format

tsən²¹⁴ mɤ kei²¹⁴ t^ha⁵⁵ pei⁵⁵ ʂɑŋ⁵¹ lai³⁵

t^ha⁵⁵ a⁵⁵ i³⁵ kei²¹⁴ tɤ t^ha⁵⁵ kei²¹⁴ tɤ

xən²¹⁴ tɕien²¹⁴ tan⁵⁵

- Moved tones from end of syllable to preceding vowel

ts ə²¹⁴ n m ɤ k e i²¹⁴ t^h a⁵⁵ p e i⁵⁵ ʂ a⁵¹ ŋ l a i³⁵

t^h a⁵⁵ a⁵⁵ i³⁵ k e i²¹⁴ t ɤ t^h a⁵⁵ k e i²¹⁴ t ɤ

x ə²¹⁴ n tɕ i e²¹⁴ n t a⁵⁵ n

- (Optionally delete tones)

Unigram word segmentation adaptor grammar

Words \rightarrow Words Word

Words \rightarrow Word

Word \rightarrow Phons

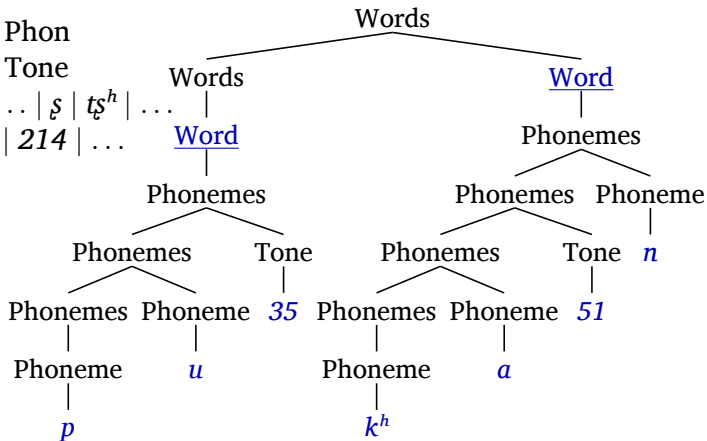
Phons \rightarrow Phon

Phons \rightarrow Phons Phon

Phons \rightarrow Phons Tone

Phon $\rightarrow ai \mid o \mid \dots \mid \xi \mid t_s^h \mid \dots$

Tone $\rightarrow 35 \mid 55 \mid 214 \mid \dots$



Collocation adaptor grammars

- Adaptor grammars with one level of collocation:

Collocs \rightarrow Colloc⁺ Colloc \rightarrow Words Words \rightarrow Word⁺

- Adaptor grammars with two levels of collocation:

Colloc2s \rightarrow Colloc2⁺ Colloc2 \rightarrow Collocs⁺
Collocs \rightarrow Colloc⁺ Colloc \rightarrow Words Words \rightarrow Word⁺

- We experiment with up to three levels of collocation here

Syllable structure adaptor grammars

- No distinction between word-internal and word-peripheral syllables

Word → Syll

Word → Syll Syll Syll

Syll → (Onset)? Rhy

Rhy → Nucleus (Coda)?

Coda → C⁺

V → ai | o | ...

Word → Syll Syll

Word → Syll Syll Syll Syll

Onset → C⁺

Nucleus → V (V | Tone)*

C → ξ | $t\xi^h$ | ...

- Distinguishing word-internal and word-peripheral syllables

Word → SyllIF

Word → SyllI Syll SyllF

SyllIF → (OnsetI)? RhyF

SyllF → (OnsetI)? RhyF

OnsetI → C⁺

CodaF → C⁺

Word → SyllI SyllF

Word → SyllI Syll Syll SyllF

SyllI → (OnsetI)? Rhy

Syll → (Onset)? Rhy

RhyF → Nucleus (CodaF)?

Mandarin Chinese word segmentation results

- Word segmentation accuracy when input *contains tones*

	Syllables		
	None	General	Specialised
Unigram	0.57	0.50	0.50
Colloc	0.69	0.67	0.67
Colloc ²	0.72	0.75	0.75
Colloc ³	0.64	0.77	0.77

- Word segmentation accuracy when *tones are removed* from input

	Syllables		
	None	General	Specialised
Unigram	0.56	0.46	0.46
Colloc	0.70	0.65	0.65
Colloc ²	0.74	0.74	0.73
Colloc ³	0.75	0.76	0.77

Comparable English results

- English word segmentation results

	Syllables		
	None	General	Specialised
Unigram	0.56	0.46	0.46
Colloc	0.74	0.67	0.66
Colloc ²	0.79	0.84	0.84
Colloc ³	0.74	0.82	0.87

Discussion of Mandarin Chinese word segmentation results

- Mandarin Chinese word segmentation results broadly consistent with English results
 - ▶ unigram segmentation accuracies are similar
 - ▶ results for other models are lower than corresponding English results
- General improvement in accuracy as number of collocation levels increases
- Caveats: the English and Mandarin Chinese corpora are not directly comparable
 - ▶ Discourse context for Mandarin Chinese corpus was far more diverse than for English corpus
 - ▶ Mandarin Chinese children were older than English children

Syllable structure and word segmentation

- Syllable structure and phonotactic constraints are very useful for English word segmentation, but are much less useful in Mandarin Chinese
 - ▶ perhaps surprising, because Mandarin Chinese has a very regular syllable structure
 - ▶ but perhaps this very predictability makes it less useful for identifying words?
 - ▶ not surprising that distinguishing word-peripheral syllables does not help, as Mandarin Chinese does not distinguish these

Tone and word segmentation

- Tones only have a small impact on segmentation accuracy
 - ▶ surprising, as they are required for lexical disambiguation
 - ▶ tones make a small improvement to simpler models (Unigram, Colloc) but no improvement with the more complex ones
 - perhaps tone is redundant given the inter-word context modelled by the Colloc²⁻³ grammars?
- *Perhaps there's a better way to represent tones in the input, or use tones in the model?*
 - ▶ Neutral tones more common on function words — perhaps this can improve segmentation accuracy?
 - ▶ Tone sandhi may give information about phonological word boundaries

Conclusion and future work

- The adaptor grammar approach to word segmentation generalises to Mandarin Chinese
- Modelling inter-word dependencies (collocations) greatly improves word segmentation accuracy in Mandarin Chinese (as in English)
- Modelling syllable structure improves segmentation accuracy by a smaller amount in Mandarin Chinese (compared to English)
- Modelling tones improves segmentation accuracy of simpler models, but not of more complex models
- Future work:
 - ▶ Comparable multi-lingual corpora of infant-directed speech
 - ▶ More realistic, richer corpora (including multi-stratal input representations)
 - ▶ Model context-sensitive dependencies (e.g., phonological rules)

Interested in computational linguistics or its applications?

We're recruiting *PhD students!*.

Contact Mark.Johnson@mq.edu.au or Katherine.Demuth@mq.edu.au for more information.

