# Adaptor Grammars:
# A framework for Bayesian
# non-parametric grammatical inference

Mark Johnson

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

Macquarie University
Sydney, Australia

MLSS "Summer School"
Updated slides available from
http://web.science.mq.edu.au/˜mjohnson/Talks.htm

MACQUARIE
UNIVERSITY

# The drunk under the lamppost

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

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

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

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

# Ideas behind talk

- Statistical methods have revolutionized computational linguistics and cognitive science
- But most successful learning methods are *parametric*
  - learn values of a *fixed number of parameters*
- *Non-parametric Bayesian methods* learn the parameters
- *Adaptor Grammars* learn probability of each *adapted subtree*
  - c.f., data-oriented parsing
- *"Rich get richer"* learning rule $\Rightarrow$ *Zipf distributions*
- Applications of Adaptor Grammars:
  - acquisition of *concatenative morphology*
  - *word segmentation* and lexical acquisition
  - topic models and *learning the referents of words*
  - learning collocations in *LDA topic models*
- Sampling (instead of EM) is a natural approach to Adaptor Grammar inference

# Language acquisition as Bayesian inference

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

- Likelihood measures how well grammar describes data
- Prior expresses knowledge of grammar before data is seen
  - can be very specific (e.g., Universal Grammar)
  - can be very general (e.g., prefer shorter grammars)
- Posterior is a *distribution* over grammars
  - captures *learner's uncertainty* about which grammar is correct
- Language learning is *non-parametric* inference
  - no (obvious) bound on number of words, grammatical morphemes, etc

MACQUARIE
UNIVERSITY

# Outline

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
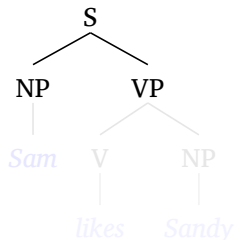- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S $\rightarrow$ NP VP |
| 0.7 | NP $\rightarrow$ *Sam* |
| 0.3 | NP $\rightarrow$ *Sandy* |
| 1 | VP $\rightarrow$ V NP |
| 0.8 | V $\rightarrow$ *likes* |
| 0.2 | V $\rightarrow$ *hates* |

S
NP   VP
Sam   V   NP
likes   Sandy

$$P(\text{Tree}) = 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
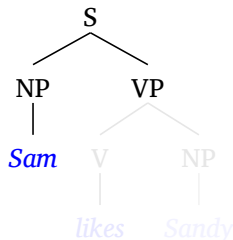- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S → NP VP |
| 0.7 | NP → *Sam* |
| 0.3 | NP → *Sandy* |
| 1 | VP → V NP |
| 0.8 | V → *likes* |
| 0.2 | V → *hates* |

$$\text{P(Tree)} = 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S $\rightarrow$ NP VP |
| 0.7 | NP $\rightarrow$ *Sam* |
| 0.3 | NP $\rightarrow$ *Sandy* |
| 1 | VP $\rightarrow$ V NP |
| 0.8 | V $\rightarrow$ *likes* |
| 0.2 | V $\rightarrow$ *hates* |

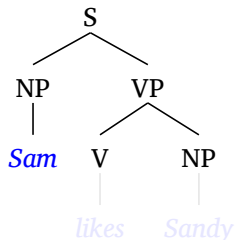$$\text{P(Tree)} = 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
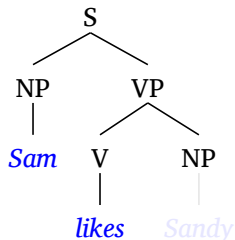- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S $\rightarrow$ NP VP |
| 0.7 | NP $\rightarrow$ *Sam* |
| 0.3 | NP $\rightarrow$ *Sandy* |
| 1 | VP $\rightarrow$ V NP |
| 0.8 | V $\rightarrow$ *likes* |
| 0.2 | V $\rightarrow$ *hates* |

$$P(\text{Tree}) = 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
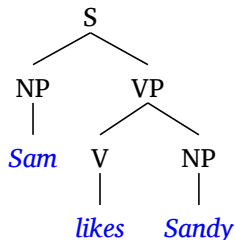- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S $\to$ NP VP |
| 0.7 | NP $\to$ *Sam* |
| 0.3 | NP $\to$ *Sandy* |
| 1 | VP $\to$ V NP |
| 0.8 | V $\to$ *likes* |
| 0.2 | V $\to$ *hates* |

$$\mathrm{P}(\text{Tree}) \;=\; 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Probabilistic context-free grammars

- Probabilistic context-free grammars (PCFGs) define *probability distributions over trees*
- Each *nonterminal node* expands by
  - choosing a rule expanding that nonterminal, and
  - recursively expanding any nonterminal children it contains
- Probability of tree is *product of probabilities of rules* used to construct it

| Probability $\theta_r$ | Rule $r$ |
|---|---|
| 1 | S $\rightarrow$ NP VP |
| 0.7 | NP $\rightarrow$ *Sam* |
| 0.3 | NP $\rightarrow$ *Sandy* |
| 1 | VP $\rightarrow$ V NP |
| 0.8 | V $\rightarrow$ *likes* |
| 0.2 | V $\rightarrow$ *hates* |

$$P(\text{Tree}) = 1 \times 0.7 \times 1 \times 0.8 \times 0.3$$

# Learning syntactic structure is hard

- Bayesian PCFG estimation works well on toy data
- Results are disappointing on "real" data
  - wrong data?
  - wrong rules?
    - *rules in PCFG must be given a priori can we learn them too?*

- Strategy: study simpler cases
  - Morphological segmentation (e.g., *walking = walk+ing*)
  - Word segmentation of unsegmented utterances
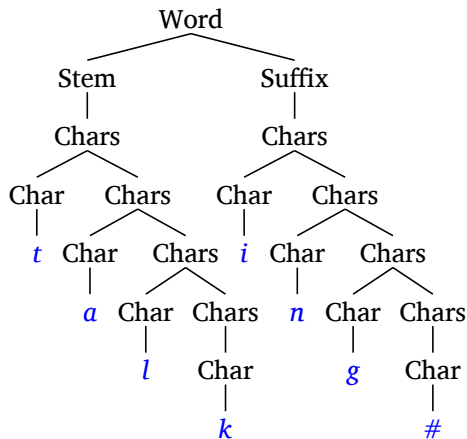
# A CFG for stem-suffix morphology

| | | | | |
|---|---|---|---|---|
| Word | → | Stem Suffix | Chars | → Char |
| Stem | → | Chars | Chars | → Char Chars |
| Suffix | → | Chars | Char | → a | b | c | ... |



- Grammar's trees can represent any segmentation of words into stems and suffixes
- ⇒ Can *represent* true segmentation
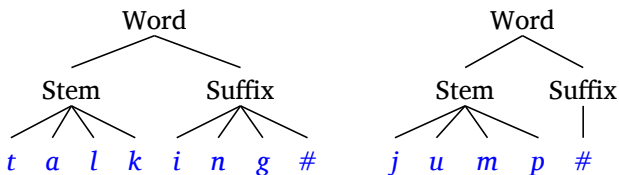- But grammar's *units of generalization (PCFG rules) are "too small"* to learn morphemes

MACQUARIE
UNIVERSITY

# A "CFG" with one rule per possible morpheme

$$
\begin{array}{rcl}
\text{Word} & \rightarrow & \text{Stem Suffix} \\
\text{Stem} & \rightarrow & \textit{all possible stems} \\
\text{Suffix} & \rightarrow & \textit{all possible suffixes}
\end{array}
$$



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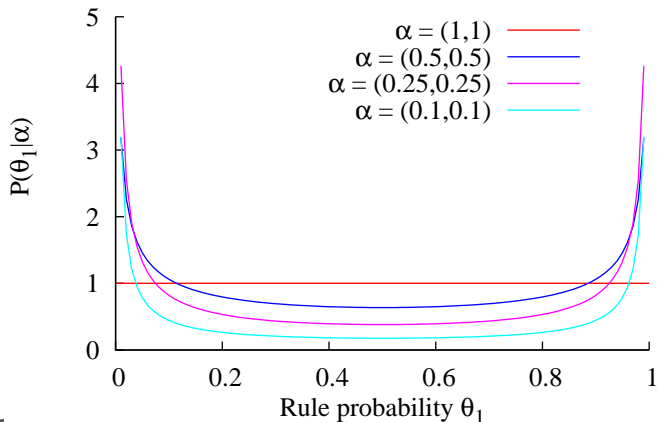  ▸ not a practical problem, as only a finite set of rules could possibly be used in any particular data set

# Maximum likelihood estimate for $\boldsymbol{\theta}$ is trivial

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

```
                          Word
                   ┌───────┴───────┐
                 Stem            Suffix
            ┌──┬──┼──┬──┐          │
            t  a  l  k  i  n  g    #
```

# Forcing generalization via sparse 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 $\theta_{A \to \beta} \approx 0$
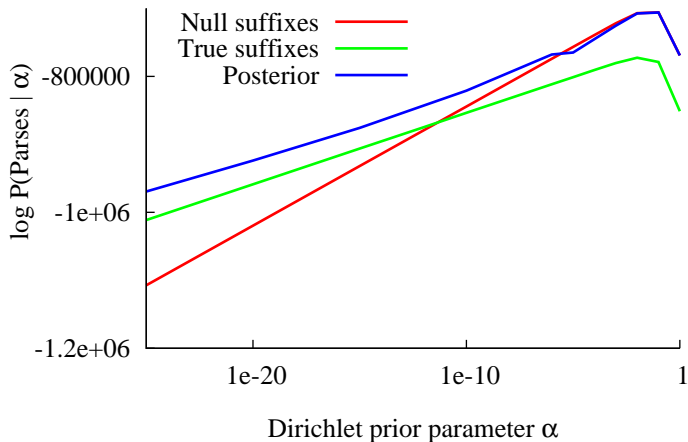- Dirichlet prior with $\alpha_{A \to \beta} \approx 0$ prefers $\theta_{A \to \beta} \approx 0$

# Morphological segmentation experiment

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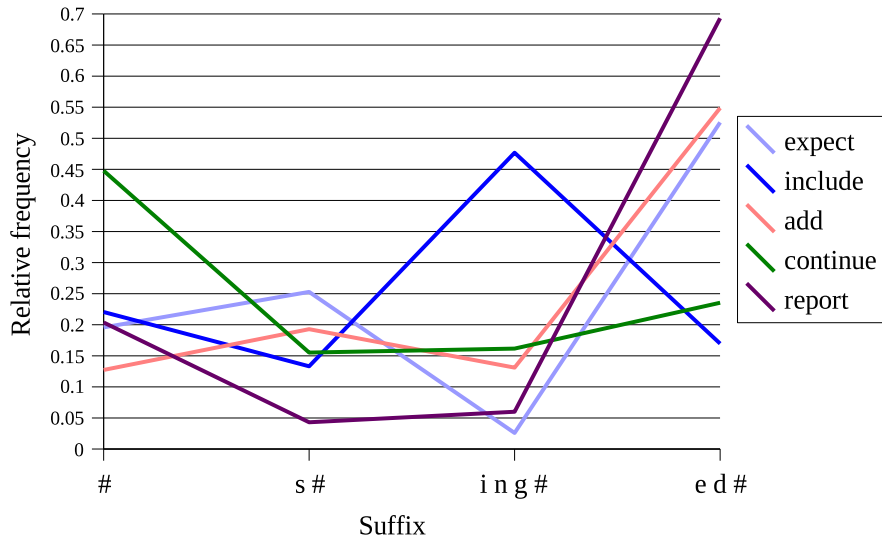
# Posterior samples from WSJ verb tokens

| $\alpha = 0.1$ | $\alpha = 10^{-5}$ | | $\alpha = 10^{-10}$ | | $\alpha = 10^{-15}$ | |
|---|---|---|---|---|---|---|
| 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

$$
\begin{aligned}
\text{Data} &= \text{``the cat chased the other cat''} \\
\text{Tokens} &= \text{``the'', ``cat'', ``chased'', ``the'', ``other'', ``cat''} \\
\text{Types} &= \text{``the'', ``cat'', ``chased'', ``other''}
\end{aligned}
$$

- Estimating $\theta$ from *word types* rather than word tokens eliminates (most) frequency variation
  - 4 common verb suffixes, so when estimating from verb types $\theta_{\text{Suffix} \to \text{i n g} \#} \approx 0.25$
- Several psycholinguists believe that humans learn morphology from word types
- Adaptor grammar mimics Goldwater et al "Interpolating between Types and Tokens" morphology-learning model

MACQUARIE
UNIVERSITY

16/117

# 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 |

MACQUARIE
UNIVERSITY

17/117

# Log posterior of models on type data



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

MACQUARIE
UNIVERSITY

# 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 naturally emerges*

MACQUARIE UNIVERSITY

# Outline

MACQUARIE
UNIVERSITY

# Bayesian inference for Dirichlet-multinomials

- Probability of next event with *uniform Dirichlet prior* with mass $\alpha$ over $m$ outcomes and observed data $\boldsymbol{Z}_{1:n} = (Z_1, \ldots, Z_n)$

$$\mathrm{P}(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}, \alpha) \quad \propto \quad n_k(\boldsymbol{Z}_{1:n}) + \alpha/m$$

where $n_k(\boldsymbol{Z}_{1:n})$ is number of times $k$ appears in $\boldsymbol{Z}_{1:n}$

- Example: Coin ($m = 2$), $\alpha = 1$, $\boldsymbol{Z}_{1:2} = (\text{heads}, \text{heads})$
  - $\mathrm{P}(Z_3 = \text{heads} \mid \boldsymbol{Z}_{1:2}, \alpha) \propto 2.5$
  - $\mathrm{P}(Z_3 = \text{tails} \mid \boldsymbol{Z}_{1:2}, \alpha) \propto 0.5$

MACQUARIE
UNIVERSITY

# Dirichlet-multinomials with many outcomes

- Predictive probability:

$$P(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}, \alpha) \;\; \propto \;\; n_k(\boldsymbol{Z}_{1:n}) + \alpha/m$$

- Suppose the number of outcomes $m \gg n$. Then:

$$P(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}, \alpha) \;\; \propto \;\; \begin{cases} n_k(\boldsymbol{Z}_{1:n}) & \text{if } n_k(\boldsymbol{Z}_{1:n}) > 0 \\[2mm] \alpha/m & \text{if } n_k(\boldsymbol{Z}_{1:n}) = 0 \end{cases}$$

- But *most outcomes will be unobserved*, so:

$$P(Z_{n+1} \notin \boldsymbol{Z}_{1:n} \mid \boldsymbol{Z}_{1:n}, \alpha) \;\; \propto \;\; \alpha$$

# From Dirichlet-multinomials to Chinese Restaurant Processes



. . .

- Suppose *number of outcomes is unbounded* but *we* pick the event labels
- If we number event types in order of occurrence
  ⇒ *Chinese Restaurant Process*

$$Z_1 = 1$$

$$P(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}, \alpha) \propto \left\{ \begin{array}{ll} n_k(\boldsymbol{Z}_{1:n}) & \text{if } k \leq m = \max(\boldsymbol{Z}_{1:n}) \\ \alpha & \text{if } k = m+1 \end{array} \right.$$
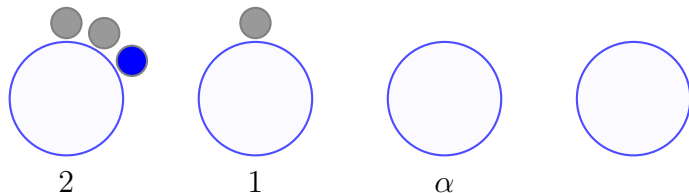
# Chinese Restaurant Process (0)



- Customer $\rightarrow$ table mapping $\boldsymbol{Z} =$
- $\mathrm{P}(\boldsymbol{z}) = 1$

- Next customer chooses a table according to:

$$\mathrm{P}(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}) \;\; \propto \;\; \left\{ \begin{array}{ll} n_k(\boldsymbol{Z}_{1:n}) & \text{if } k \leq m = \max(\boldsymbol{Z}_{1:n}) \\ \alpha & \text{if } k = m + 1 \end{array} \right.$$

# Chinese Restaurant Process (1)



- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1$
- $P(\boldsymbol{z}) = \alpha/\alpha$

- Next customer chooses a table according to:

$$P(Z_{n+1} = k \mid \boldsymbol{Z}_{1:n}) \quad \propto \quad \left\{ \begin{array}{ll} n_k(\boldsymbol{Z}_{1:n}) & \text{if } k \leq m = \max(\boldsymbol{Z}_{1:n}) \\ \alpha & \text{if } k = m+1 \end{array} \right.$$

# Chinese Restaurant Process (2)



- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1, 1$
- $P(\boldsymbol{z}) = \alpha/\alpha \times 1/(1 + \alpha)$

- Next customer chooses a table according to:

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

# Chinese Restaurant Process (3)



- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1, 1, 2$
- $P(\boldsymbol{z}) = \alpha/\alpha \times 1/(1+\alpha) \times \alpha/(2+\alpha)$

- Next customer chooses a table according to:

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

# Chinese Restaurant Process (4)



- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1, 1, 2, 1$
- $P(\boldsymbol{z}) = \alpha/\alpha \times 1/(1+\alpha) \times \alpha/(2+\alpha) \times 2/(3+\alpha)$

- Next customer chooses a table according to:

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

# Labeled Chinese Restaurant Process (0)

- Table $\rightarrow$ label mapping $\boldsymbol{Y} =$
- Customer $\rightarrow$ table mapping $\boldsymbol{Z} =$
- Output sequence $\boldsymbol{X} =$
- $\mathrm{P}(\boldsymbol{X}) = 1$

- *Base distribution* $\mathrm{P}_0(Y)$ generates a *label* $Y_k$ for each table $k$
- All customers sitting at table $k$ (i.e., $Z_i = k$) share label $Y_k$
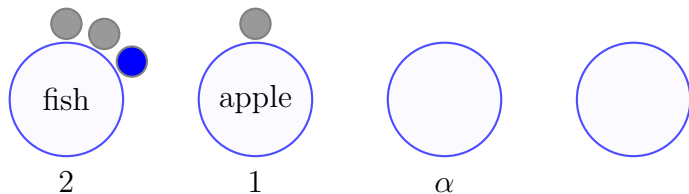- Customer $i$ sitting at table $Z_i$ has label $X_i = Y_{Z_i}$

# Labeled Chinese Restaurant Process (1)



- Table → label mapping $\boldsymbol{Y}$ = fish
- Customer → table mapping $\boldsymbol{Z}$ = 1
- Output sequence $\boldsymbol{X}$ = fish
- $\mathrm{P}(\boldsymbol{X}) = \alpha/\alpha \times \mathrm{P}_0(\text{fish})$

- *Base distribution* $\mathrm{P}_0(Y)$ 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 $\boldsymbol{Y} =$ fish
- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1, 1$
- Output sequence $\boldsymbol{X} =$ fish,fish
- $P(\boldsymbol{X}) = P_0(\text{fish}) \times 1/(1 + \alpha)$

- *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 $\boldsymbol{Y} =$ fish,apple
- Customer $\rightarrow$ table mapping $\boldsymbol{Z} = 1, 1, 2$
- Output sequence $\boldsymbol{X} =$ fish,fish,apple
- $P(\boldsymbol{X}) = P_0(\text{fish}) \times 1/(1 + \alpha) \times \alpha/(2 + \alpha)P_0(\text{apple})$

- *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 → label mapping $\boldsymbol{Y}$ = fish,apple
- Customer → table mapping $\boldsymbol{Z}$ = 1, 1, 2
- Output sequence $\boldsymbol{X}$ = fish,fish,apple,fish
- $P(\boldsymbol{X}) = P_0(\text{fish}) \times 1/(1+\alpha) \times \alpha/(2+\alpha) P_0(\text{apple}) \times 2/(3+\alpha)$

- *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) generalise 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) generalise CRPs with an additional concentration parameter
  - this parameter specifies the asymptotic power-law behaviour

- *Labeled CRPs* use a *base distribution* to define distributions over arbitrary objects
  - base distribution *"labels the tables"*
  - base distribution can have *infinite support*
  - concentrates mass on a countable subset
  - power-law behaviour $\Rightarrow$ Zipfian distributions

# Nonparametric extensions of PCFGs

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

MACQUARIE
UNIVERSITY

# Outline

MACQUARIE
UNIVERSITY

# 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

# 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* $B \to \beta \in R$ with prob. $\theta_{B \to \beta}$, and recursively expand nonterminals in $\beta$
  - ▸ 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* $B \to \beta \in R$ with prob. $\propto \alpha_B \theta_{B \to \beta}$, and recursively expand nonterminals in $\beta$

# Adaptor grammar for stem-suffix morphology



Word → Stem Suffix
Stem → Phons
Suffix → Phons
Phons → Phon
Phons → Phon Phons

or in *abbreviated form* with
non-adapted nonterminals suppressed

Word → Stem Suffix
Stem → Phon$^+$
Suffix → Phon$^+$

# Adaptor grammar for stem-suffix morphology (0)

<u>Word</u> → Stem Suffix

<u>Stem</u> → Phoneme$^+$

<u>Suffix</u> → Phoneme$^\star$

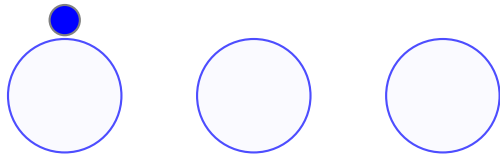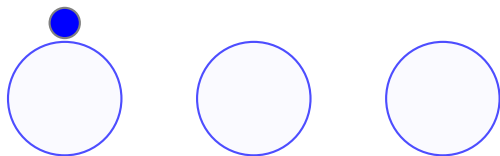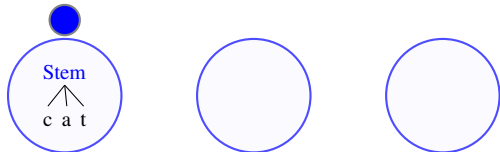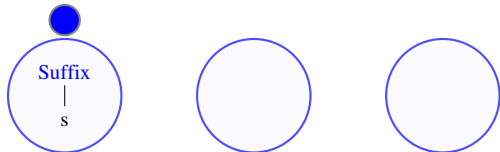Generated words:

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

<u>Word</u> → Stem Suffix

<u>Stem</u> → Phoneme$^+$

<u>Suffix</u> → Phoneme$^\star$

Generated words:

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

<u>Word</u> → Stem Suffix

<u>Stem</u> → Phoneme⁺

<u>Suffix</u> → Phoneme⋆

Generated words:

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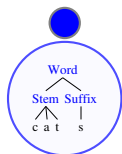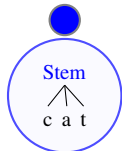


Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$

Generated words:

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



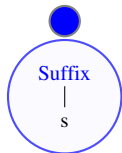Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$

Generated words: cats

# Adaptor grammar for stem-suffix morphology (2a)
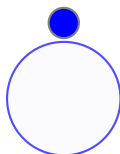
Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$



Generated words: cats

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

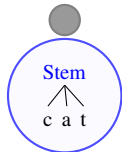

Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$

Generated words: cats

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



<u>Word</u> → Stem Suffix

<u>Stem</u> → Phoneme$^+$

<u>Suffix</u> → Phoneme$^\star$

Generated words: cats

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



Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$
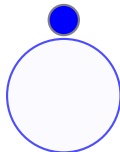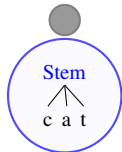
Generated words: cats, dogs

# Adaptor grammar for stem-suffix morphology (3)



Word → Stem Suffix

Stem → Phoneme$^+$

Suffix → Phoneme$^\star$

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 $A \to \beta$ with probability $\theta_{A \to \beta}$
- Each adapted nonterminal $A$ is associated with a CRP (or PYP) that caches previously generated subtrees rooted in $A$
- An *adapted nonterminal* $A$ expands:
  - to a subtree $\tau$ rooted in $A$ with probability proportional to the number of times $\tau$ was previously generated
  - using $A \to \beta$ with probability proportional to $\alpha_A \theta_{A \to \beta}$

MACQUARIE
UNIVERSITY

# 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* $A \to \beta$, where $A \in N$ and $\beta \in (N \cup W)^\star$
- 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:

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

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

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

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 $\boldsymbol{\theta}$, where:

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

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

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

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

$$
\mathrm{TD}_A(G_1, \dots, G_n) \left( \underbrace{\phantom{t_1 \dots t_n}}_{\substack{A \\ t_1 \dots t_n}} \right) = \prod_{i=1}^{n} G_i(t_i).
$$

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

MACQUARIE
UNIVERSITY

# 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 $A \in N$, $\alpha_A$ is the parameter of the Dirichlet process associated with $A$.

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

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

The grammar generates the distribution $G_S$.
One Dirichlet Process for each adapted non-terminal $A$ (i.e., $\alpha_A > 0$).

MACQUARIE
UNIVERSITY

# 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., $\alpha_B$) can themselves be estimated (e.g., slice sampling)

# 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

MACQUARIE
UNIVERSITY

# Unsupervised word segmentation

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

$$j \vartriangle u \blacktriangle w \vartriangle a \vartriangle n \vartriangle t \blacktriangle t \vartriangle u \blacktriangle s \vartriangle i \blacktriangle ð \vartriangle ə \blacktriangle b \vartriangle ʊ \vartriangle k$$
"you want to see the book"

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

# CFG models of word segmentation

Words → Word
Words → Word Words
Word → Phons
Phons → Phon
Phons → Phon Phons
Phon → $a \mid b \mid$ ...

- 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 *bʊk* is a Word

# Towards non-parametric grammars

Words → Word
Words → Word Words
Word → *all possible phoneme sequences*

- Learn probability Word → *b ʊ k*
- But *infinitely many possible Word expansions*
  ⇒ this grammar is *not a PCFG*

- Given *fixed training data*, only finitely many useful rules
  ⇒ use data to choose Word rules as well as their probabilities

- An Adaptor Grammar can do precisely this!

# Unigram adaptor grammar (Brent)

Words → Word
Words → Word Words
<u>Word</u> → Phons
Phons → Phon
Phons → Phon Phons

- <u>Word</u> nonterminal is adapted
⇒ To generate a <u>Word</u>:
    ‣ select a previously generated <u>Word</u> subtree
      with prob. ∝ number of times it has been generated
    ‣ expand using <u>Word</u> → Phons rule with prob. ∝ $\alpha_{\text{Word}}$
      and recursively expand Phons

# 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 → Word⁺
<u>Word</u> → Phoneme⁺



- 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 → <u>Word</u> Words | 1 | Words → <u>Word</u> |
  |---|---|---|---|
  | 1 | <u>Word</u> → Phon | | |
  | 1 | Phons → Phon Phons | 1 | Phons → Phon |
  | 1 | Phon → $D$ | 1 | Phon → $G$ |
  | 1 | Phon → $A$ | 1 | Phon → $E$ |

- **A grammar learnt from Brent corpus**

  | 16625 | Words → <u>Word</u> Words | 9791 | Words → <u>Word</u> |
  |---|---|---|---|
  | 1575 | <u>Word</u> → Phons | | |
  | 4962 | Phons → Phon Phons | 1575 | Phons → Phon |
  | 134 | Phon → $D$ | 41 | Phon → $G$ |
  | 180 | Phon → $A$ | 152 | Phon → $E$ |
  | 460 | <u>Word</u> → (Phons (Phon $y$) (Phons (Phon $u$))) | | |
  | 446 | <u>Word</u> → (Phons (Phon $w$) (Phons (Phon $A$) (Phons (Phon $t$) | | |
  | 374 | <u>Word</u> → (Phons (Phon $D$) (Phons (Phon $6$))) | | |
  | 372 | <u>Word</u> → (Phons (Phon $&$) (Phons (Phon $n$) (Phons (Phon $d$) | | |

# Undersegmentation errors with Unigram model

$$\text{Words} \rightarrow \underline{\text{Word}}^+ \qquad \underline{\text{Word}} \rightarrow \text{Phon}^+$$

- 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 ⇒ Words

$$\text{Sentence} \to \text{Colloc}^+$$
$$\underline{\text{Colloc}} \to \text{Word}^+$$
$$\underline{\text{Word}} \to \text{Phon}^+$$

```
                        Sentence
        _____|_____
       |                |            |
     Colloc           Colloc       Colloc
     __|____            |          __|____
    |       |           |         |       |
  Word    Word        Word      Word     Word
  /|\    /|||\\        /|\       /\      /|\
 j u   w a n t t     u s i      ð ə     b ʊ k
```

- 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; ≈ Goldwater's bigram model)

# Collocations $\Rightarrow$ Words $\Rightarrow$ Syllables

$$\text{Sentence} \rightarrow \text{Colloc}^+ \qquad \underline{\text{Colloc}} \rightarrow \text{Word}^+$$

$$\underline{\text{Word}} \rightarrow \text{Syllable}^{\{1:3\}} \qquad \text{Syllable} \rightarrow (\text{Onset})\ \text{Rhyme}$$

$$\underline{\text{Onset}} \rightarrow \text{Consonant}^+ \qquad \text{Rhyme} \rightarrow \text{Nucleus}\ (\text{Coda})$$

$$\underline{\text{Nucleus}} \rightarrow \text{Vowel}^+ \qquad \underline{\text{Coda}} \rightarrow \text{Consonant}^+$$

Sentence

Colloc        Colloc

Word     Word     Word

Onset Nucleus Coda Nucleus Coda Onset Nucleus Coda

l     ʊ     k     æ     t     ð     ɪ     s

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

# Distinguishing internal onsets/codas helps

Sentence → Colloc$^+$

Word → SyllableIF

Word → SyllableI Syllable SyllableF

OnsetI → Consonant$^+$

Nucleus → Vowel$^+$

Colloc → Word$^+$

Word → SyllableI SyllableF

SyllableIF → (OnsetI) RhymeI

RhymeF → Nucleus (CodaF)

CodaF → Consonant$^+$



- With 2 Collocation levels, not distinguishing initial/final clusters, f-score = 84%
- With 3 Collocation levels, distinguishing initial/final clusters, f-score = 87%

# Collocations$^2$ ⇒ Words ⇒ Syllables

# Summary of English word segmentation

- 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

# Outline

MACQUARIE
UNIVERSITY

# 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 $PRMPT
  - corpus contains 50,118 utterances, 187,533 word tokens
    zen3me gei3 ta1 bei1 shang4 lai2 (1.) ?
    ta1: (.) a1yi2 gei3 de (.) ta1 gei3 de .
    hen3 jian3dan1 .
- Used Pinyin to IPA translation program to produce IPA:
    tsən$^{214}$mɤ kei$^{214}$ tʰa$^{55}$ pei$^{55}$ ʂɑŋ$^{51}$ lai$^{35}$
    tʰa$^{55}$ a$^{55}$i$^{35}$ kei$^{214}$ tɤ tʰa$^{55}$ kei$^{214}$ tɤ
    xən$^{214}$ tɕiɛn$^{214}$tan$^{55}$
- Moved tones from end of syllable to preceding vowel
    ts ə $^{214}$ n m ɤ k e i $^{214}$ tʰ a $^{55}$ p e i $^{55}$ ʂ ɑ $^{51}$ ŋ l ai $^{35}$
    tʰ a $^{55}$ a $^{55}$ i $^{35}$ k e i $^{214}$ t ɤ tʰ a $^{55}$ k e i $^{214}$ t ɤ
    x ə $^{214}$ n tɕ iɛ $^{214}$ n t a $^{55}$ n
- (Optionally delete tones)

MACQUARIE
UNIVERSITY

# Unigram word segmentation adaptor grammar

Words → Words <u>Word</u>
Words → <u>Word</u>
<u>Word</u> → Phons
Phons → Phon
Phons → Phons Phon
Phons → Phons Tone
Phon → $ai \mid t \mid \ldots$
Tone → $35 \mid 55 \mid 214 \mid \ldots$

# Collocation adaptor grammars

- Adaptor grammars with one level of collocation:

$$\text{Collocs} \rightarrow \underline{\text{Colloc}}^+ \qquad \underline{\text{Colloc}} \rightarrow \text{Words}$$
$$\text{Words} \rightarrow \underline{\text{Word}}^+$$

- Adaptor grammars with two levels of collocation:

$$\text{Colloc2s} \rightarrow \underline{\text{Colloc2}}^+ \qquad \underline{\text{Colloc2}} \rightarrow \text{Collocs}^+$$
$$\text{Collocs} \rightarrow \underline{\text{Colloc}}^+ \qquad \underline{\text{Colloc}} \rightarrow \text{Words}$$
$$\text{Words} \rightarrow \underline{\text{Word}}^+$$

- We experiment with *up to three collocation levels* here

# Syllable structure adaptor grammars

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

  $$\underline{\text{Word}} \rightarrow \text{Syll} \qquad\qquad \underline{\text{Word}} \rightarrow \text{Syll Syll}$$
  $$\underline{\text{Word}} \rightarrow \text{Syll Syll Syll} \qquad \underline{\text{Word}} \rightarrow \text{Syll Syll Syll Syll}$$
  $$\text{Syll} \rightarrow (\underline{\text{Onset}})^{?} \ \underline{\text{Rhy}} \qquad \underline{\text{Onset}} \rightarrow \text{C}^{+}$$
  $$\underline{\text{Rhy}} \rightarrow \underline{\text{Nucleus}} \ (\underline{\text{Coda}})^{?} \qquad \underline{\text{Nucleus}} \rightarrow \text{V } (\text{V } | \text{ Tone})^{\star}$$
  $$\underline{\text{Coda}} \rightarrow \text{C}^{+} \qquad\qquad \text{C} \rightarrow \ | \ t \ | \dots$$
  $$\text{V} \rightarrow ai \ | \ o \ | \dots$$

- *Distinguishing word-internal and word-peripheral syllables*

  $$\underline{\text{Word}} \rightarrow \text{SyllIF} \qquad\qquad \underline{\text{Word}} \rightarrow \text{SyllI SyllF}$$
  $$\underline{\text{Word}} \rightarrow \text{SyllI Syll SyllF} \qquad \underline{\text{Word}} \rightarrow \text{SyllI Syll Syll SyllF}$$
  $$\text{SyllIF} \rightarrow (\underline{\text{OnsetI}})^{?} \ \underline{\text{RhyF}} \qquad \text{SyllI} \rightarrow (\underline{\text{OnsetI}})^{?} \ \underline{\text{Rhy}}$$
  $$\text{SyllF} \rightarrow (\underline{\text{OnsetI}})^{?} \ \underline{\text{RhyF}} \qquad \text{Syll} \rightarrow (\underline{\text{Onset}})^{?} \ \underline{\text{Rhy}}$$
  $$\underline{\text{OnsetI}} \rightarrow \text{C}^{+} \qquad\qquad \underline{\text{RhyF}} \rightarrow \underline{\text{Nucleus}} \ (\underline{\text{CodaF}})^{?}$$
  $$\underline{\text{CodaF}} \rightarrow \text{C}^{+}$$

MACQUARIE
UNIVERSITY

# 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        |
| $\text{Colloc}^2$    | 0.72      | 0.75    | 0.75        |
| $\text{Colloc}^3$    | *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        |
| $\text{Colloc}^2$    | 0.74      | 0.74    | 0.73        |
| $\text{Colloc}^3$    | 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 similiar
  - 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

MACQUARIE
UNIVERSITY

# 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 $\text{Colloc}^{2-3}$ 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

# Outline

# Two hypotheses about language acquisition

1. Pre-programmed *staged acquisition* of linguistic components
   - Conventional view of *lexical acquisition*, e.g., Kuhl (2004)
     - child first learns the phoneme inventory, which it then uses to learn
     - phonotactic cues for word segmentation, which are used to learn
     - phonological forms of words in the lexicon, . . .

2. *Interactive acquisition* of all linguistic components together
   - corresponds to *joint inference* for all components of language
   - stages in language acquisition might be due to:
     - child's input may contain more information about some components
     - some components of language may be learnable with less data

# Synergies: an advantage of interactive learning

- An *interactive learner* can take advantage of *synergies in acquisition*
  - partial knowledge of component $A$ provides information about component $B$
  - partial knowledge of component $B$ provides information about component $A$
- A staged learner can only take advantage of one of these dependencies
- An interactive or *joint learner* can benefit from a positive feedback cycle between $A$ and $B$
- Are there synergies in *learning how to segment words* and *learning the referents of words*?

MACQUARIE
UNIVERSITY

# Prior work: mapping words to referents



- Input to learner:
  - word sequence: *Is that the pig?*
  - objects in nonlinguistic context: DOG, PIG
- Learning objectives:
  - identify utterance topic: PIG
  - identify word-topic mapping: $pig \mapsto PIG$

# Frank et al (2009) "topic models" as PCFGs

- Prefix sentences with *possible topic marker*, e.g., PIG|DOG

- PCFG rules *choose a topic* from topic marker and *propagate it through sentence*

- Each word is either generated from sentence topic or null topic $\emptyset$

```
                          Sentence
                             |
                          Topic_pig
                          /        \
                     Topic_pig    Word_pig
                     /       \        |
                Topic_pig   Word_∅   pig
                /      \       |
           Topic_pig  Word_∅  the
           /     \       |
      Topic_pig Word_∅  that
         |        |
      PIG|DOG    is
```

- Grammar can require *at most one topical word per sentence*

- Bayesian inference for PCFG rules and trees corresponds to Bayesian inference for word and sentence topics using topic model (Johnson 2010)

# Word segmentation with adaptor grammars

- Adaptor grammars (AGs) can learn the probability of entire subtrees (as well as rules)
- AGs can express several different word segmentation models
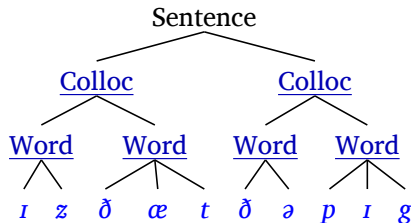- Learning collocations as well as words significantly improves segmentation accuracy

Sentence → $\underline{\text{Colloc}}^+$
$\underline{\text{Colloc}}$ → $\underline{\text{Word}}^+$
$\underline{\text{Word}}$ → Phon$^+$

# AGs for joint segmentation and referent-mapping

- Combine topic-model PCFG with word segmentation AGs
- Input consists of unsegmented phonemic forms prefixed with possible topics:

$$\text{PIG}\big|\text{DOG } \textit{ɪ z ð æ t ð ə p ɪ g}$$

- E.g., combination of *Frank "topic model"* and *unigram segmentation model*
    - equivalent to Jones et al (2010)

- Easy to define *other combinations of topic models and segmentation models*

```
Sentence
   |
Topic_pig
   /        \
Topic_pig    Word_pig
   /    \      /  |  \
Topic_pig  Word_∅  p  ɪ  g
   /    \      ð  ə
Topic_pig  Word_∅  ð  æ  t
   /    \
Topic_pig  Word_∅
   |        ð  æ  t
PIG|DOG  ɪ  z
```

MACQUARIE
UNIVERSITY

# Collocation topic model AG



- Collocations are either "topical" or not
- Easy to modify this grammar so
  - at most one topical word per sentence, or
  - at most *one topical word per topical collocation*

# Experimental set-up

- Input consists of unsegmented phonemic forms prefixed with possible topics:

  PIG|DOG *ɪ z ð æ t ð ə p ɪ g*

  - Child-directed speech corpus collected by Fernald et al (1993)
  - Objects in visual context annotated by Frank et al (2009)

- Bayesian inference for AGs using MCMC (Johnson et al 2009)
  - Uniform prior on PYP $a$ parameter
  - "Sparse" Gamma$(100, 0.01)$ on PYP $b$ parameter

- For each grammar we ran 8 MCMC chains for 5,000 iterations
  - collected word segmentation and topic assignments at every 10th iteration during last 2,500 iterations
    $\Rightarrow$ 2,000 sample analyses per sentence
  - computed and evaluated the modal (i.e., most frequent) sample analysis of each sentence

MACQUARIE UNIVERSITY

# Does non-linguistic context help segmentation?

| Model | | word segmentation |
| segmentation | topics | token f-score |
|:---:|:---:|:---:|
| unigram | not used | 0.533 |
| unigram | any number | 0.537 |
| unigram | one per sentence | 0.547 |
| collocation | not used | 0.695 |
| collocation | any number | 0.726 |
| collocation | one per sentence | 0.719 |
| collocation | one per collocation | **0.750** |

- Not much improvement with unigram model
  - consistent with results from Jones et al (2010)
- Larger improvement with collocation model
  - most gain with *one topical word per topical collocation*
    (this constraint cannot be imposed on unigram model)

# Does better segmentation help topic identification?

- Task: identify object (if any) *this sentence* is about

| Model | | sentence referent | |
|:---:|:---:|:---:|:---:|
| **segmentation** | **topics** | **accuracy** | **f-score** |
| unigram | not used | 0.709 | 0 |
| unigram | any number | 0.702 | 0.355 |
| unigram | one per sentence | 0.503 | 0.495 |
| collocation | not used | 0.709 | 0 |
| collocation | any number | 0.728 | 0.280 |
| collocation | one per sentence | 0.440 | 0.493 |
| collocation | one per collocation | **0.839** | **0.747** |

- The collocation grammar with *one topical word per topical collocation* is the only model clearly better than baseline

# Does better segmentation help learning word-to-referent mappings?

- Task: identify *head nouns* of NPs referring to topical objects
  (e.g. $pig \mapsto$ PIG in input PIG | DOG $ɪ z ð æ t ð ə p ɪ g$)

| Model | | topical word |
| segmentation | topics | f-score |
| --- | --- | --- |
| unigram | not used | 0 |
| unigram | any number | 0.149 |
| unigram | one per sentence | 0.147 |
| collocation | not used | 0 |
| collocation | any number | 0.220 |
| collocation | one per sentence | 0.321 |
| collocation | one per collocation | **0.636** |

- The collocation grammar with one topical word per topical collocation is best at identifying head nouns of referring NPs

MACQUARIE
UNIVERSITY

# Summary of segmentation and word-to-referent mappings

- *Word to object mapping is learnt more accurately when words are segmented more accurately*
  - ▸ improving segmentation accuracy improves topic detection and acquisition of topical words
- *Word segmentation accuracy improves when exploiting non-linguistic context information*
  - ▸ incorporating word-topic mapping improves segmentation accuracy (at least with collocation grammars)
- ⇒ *There seem to be synergies a learner could exploit when learning word segmentation and word-object mappings*
  - ▸ Caveat: results seem to depend on details of model
- Complexity of models limited by ability to "pass features" in a PCFG
  - ▸ future work: extend the AG framework to permit "feature-passing"

# Outline

MACQUARIE UNIVERSITY

# LDA topic models

- LDA topic models are *admixture models* of documents
  - ‣ topics are assigned to *words* (not sentences or documents)
- An LDA topic model learns:
  - ‣ the topics expressed in a document
  - ‣ the words characteristic of a topic
- Each topic $i$ is a distribution over words $\boldsymbol{\phi}_i$
- Each document $j$ has a *distribution* $\boldsymbol{\theta}_j$ over topics
- To generate document $j$:
  - ‣ for each word position in document:
    - – choose a topic $z$ according to $\boldsymbol{\theta}_j$, and then
    - – choose a word belonging to that topic according to $\boldsymbol{\phi}_z$
- "Sparse priors" on $\boldsymbol{\phi}$ and $\boldsymbol{\theta}$
  - $\Rightarrow$ most documents have few topics
  - $\Rightarrow$ most topics have few words

# LDA topic models as Bayes nets

$$
\begin{aligned}
\boldsymbol{\phi}_i &\sim \mathrm{Dir}(\boldsymbol{\beta}) & i &= 1, \ldots, \ell = \text{number of topics} \\
\boldsymbol{\theta}_j &\sim \mathrm{Dir}(\boldsymbol{\alpha}) & j &= 1, \ldots, m = \text{number of documents} \\
z_{j,k} &\sim \boldsymbol{\theta}_j & j &= 1, \ldots, m \\
& & k &= 1, \ldots, n = \text{number of words in a document} \\
w_{j,k} &\sim \boldsymbol{\phi}_{z_{j,k}} & j &= 1, \ldots, m \\
& & k &= 1, \ldots, n
\end{aligned}
$$

# LDA topic models as PCFGs (1)

- Prefix strings from document $j$ with a *document identifier* "$_{-j}$"

$$
\begin{aligned}
&\text{Sentence} \rightarrow \text{Doc}'_j && j \in 1, \ldots, m \\
&\text{Doc}'_j \rightarrow \ _{-j} && j \in 1, \ldots, m \\
&\text{Doc}'_j \rightarrow \text{Doc}'_j\,\text{Doc}_j && j \in 1, \ldots, m \\
&\text{Doc}_j \rightarrow \text{Topic}_i && i \in 1, \ldots, \ell \\
& && j \in 1, \ldots, m \\
&\text{Topic}_i \rightarrow \ w && i \in 1, \ldots, \ell \\
& && w \in \mathcal{V}
\end{aligned}
$$

# LDA topic models as PCFGs (2)

- Spine *propagates document id up through tree*

$$\text{Sentence} \rightarrow \text{Doc}'_j \qquad j \in 1, \dots, m$$
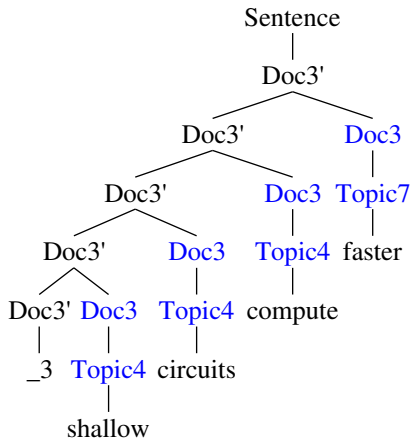$$\text{Doc}'_j \rightarrow \,\_j \qquad\qquad j \in 1, \dots, m$$
$$\text{Doc}'_j \rightarrow \text{Doc}'_j \, \text{Doc}_j \qquad j \in 1, \dots, m$$
$$\text{Doc}_j \rightarrow \text{Topic}_i \qquad i \in 1, \dots, \ell$$
$$\qquad\qquad\qquad\qquad j \in 1, \dots, m$$
$$\text{Topic}_i \rightarrow w \qquad\qquad i \in 1, \dots, \ell$$
$$\qquad\qquad\qquad\qquad w \in \mathcal{V}$$

```
                    Sentence
                       |
                    Doc3'
                   /      \
               Doc3'       Doc3
              /    \          |
          Doc3'    Doc3    Topic7
          /   \      |        |
      Doc3'  Doc3  Topic4   faster
      /  \     |     |
  Doc3' Doc3 Topic4 compute
  /  \    |    |
_3 Topic4 circuits
     |
  shallow
```

# LDA topic models as PCFGs (3)

- $\text{Doc}_j \to \text{Topic}_i$ rules map *documents to topics*

$$\begin{aligned}
\text{Sentence} &\to \text{Doc}'_j & j &\in 1, \ldots, m \\
\text{Doc}'_j &\to {}_{-j} & j &\in 1, \ldots, m \\
\text{Doc}'_j &\to \text{Doc}'_j \, \text{Doc}_j & j &\in 1, \ldots, m \\
\text{Doc}_j &\to \text{Topic}_i & i &\in 1, \ldots, \ell \\
& & j &\in 1, \ldots, m \\
\text{Topic}_i &\to w & i &\in 1, \ldots, \ell \\
& & w &\in \mathcal{V}
\end{aligned}$$

# LDA topic models as PCFGs (4)

- $\text{Topic}_i \to w$ rules map *topics to words*

$$
\begin{aligned}
\text{Sentence} &\to \text{Doc}'_j & j \in 1, \ldots, m \\
\text{Doc}'_j &\to \_{-j} & j \in 1, \ldots, m \\
\text{Doc}'_j &\to \text{Doc}'_j \, \text{Doc}_j & j \in 1, \ldots, m \\
\text{Doc}_j &\to \text{Topic}_i & i \in 1, \ldots, \ell \\
& & j \in 1, \ldots, m \\
\text{Topic}_i &\to w & i \in 1, \ldots, \ell \\
& & w \in \mathcal{V}
\end{aligned}
$$

# Topic model with collocations

- Combines *PCFG topic model* and *segmentation adaptor grammar*

$$\text{Sentence} \rightarrow \text{Doc}_j \qquad j \in 1, \ldots, m$$
$$\text{Doc}_j \rightarrow {}_{-j} \qquad j \in 1, \ldots, m$$
$$\text{Doc}_j \rightarrow \text{Doc}_j \, \text{Topic}_i \qquad i \in 1, \ldots, \ell;$$
$$\qquad\qquad\qquad\qquad\qquad j \in 1, \ldots, m$$
$$\underline{\text{Topic}_i} \rightarrow \text{Words} \qquad i \in 1, \ldots, \ell$$
$$\text{Words} \rightarrow \text{Word}$$
$$\text{Words} \rightarrow \text{Words} \, \text{Word}$$
$$\text{Word} \rightarrow w \qquad w \in \mathcal{V}$$

# Finding topical collocations in NIPS abstracts

- Run topical collocation adaptor grammar on NIPS corpus
- Run with $\ell = 20$ topics (i.e., 20 distinct $\text{Topic}_i$ nonterminals)
- Corpus is segmented by punctuation
  - ▸ terminal strings are fairly short
  - ⇒ inference is fairly efficient
- Used standard AG implementation
  - ▸ Pitman-Yor adaptors
  - ▸ sampled Pitman-Yor $a$ and $b$ parameters
  - ▸ flat and "vague Gamma" priors on Pitman-Yor $a$ and $b$ parameters

MACQUARIE
UNIVERSITY

# 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 skeletal 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 ocularity) (T_12 drift rate)

MACQUARIE
UNIVERSITY

# Outline

MACQUARIE
UNIVERSITY

# What do we have to learn?

- To learn an adaptor grammar, we need:
  - ▸ probabilities of grammar rules
  - ▸ adapted subtrees and their probabilities for adapted non-terminals
- If we knew the true parse trees for a training corpus, we could:
  - ▸ read off the adapted subtrees from the corpus
  - ▸ count rules and adapted subtrees in corpus
  - ▸ compute the rule and subtree probabilities from these counts
    - – simple computation (smoothed relative frequencies)
- If we aren't given the parse trees:
  - ▸ there can be *infinitely many* possible adapted subtrees
  - ⇒ can't track the probability of all of them (as in EM)
  - ▸ but *sample parses of a finite corpus* only include finitely many
- Sampling-based methods learn the relevant subtrees as well as their weights

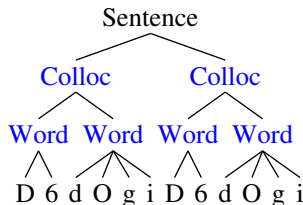# A Gibbs sampler for learning adaptor grammars

- Gibbs sampling for learning adaptor grammars
  - Assign (random) parse trees to each sentence, and compute rule and subtree counts
  - Repeat forever:
    - pick a sentence (and corresponding parse) at random
    - deduct the counts for the sentence's parse from current rule and subtree counts
    - sample a parse for sentence according to updated grammar
    - add sampled parse's counts to rule and subtree counts

- Sampled parse trees and grammar converges to Bayesian posterior distribution

# Sampling parses from an adaptor grammar

- Sampling a parse tree for a sentence is computationally most demanding part of learning algorithm
- Component-wise Metropolis-within-Gibbs sampler for parse trees:
  - adaptor grammar rules and probabilities *change on the fly*
  - construct PCFG *proposal grammar* from adaptor grammar for previous sentences
  - sample a parse from PCFG proposal grammar
  - use accept/reject to convert samples from proposal PCFG to samples from adaptor grammar
- For particular adaptor grammars, there are often more efficient algorithms

MACQUARIE
UNIVERSITY

# Details about sampling parses

- Adaptor grammars are *not context-free*
- The probability of a rule (and a subtree) can change within a single sentence
  - breaks standard dynamic programming



- But with moderate or large corpora, the probabilities don't change by much
  - use Metropolis-Hastings accept/reject with a PCFG proposal distribution
- Rules of PCFG proposal grammar $G'(\boldsymbol{t}_{-j})$ consist of:
  - rules $A \to \beta$ from base PCFG: $\theta'_{A \to \beta} \propto \alpha_A \theta_{A \to \beta}$
  - A rule $A \to \text{YIELD}(t)$ for each table $t$ in $A$'s restaurant: $\theta'_{A \to \text{YIELD}(t)} \propto n_t$, the number of customers at table $t$
- Map parses using $G'(\boldsymbol{t}_{-j})$ back to adaptor grammar parses

# Random vs incremental initialization

- The Gibbs sampler parse trees $\boldsymbol{t}$ needs to be initialized somehow

  Random initialization: Assign each string $x_i$ a random parse $t_i$ generated by base PCFG

  Incremental initialization: Sample $t_i$ from $P(t \mid x_i, \boldsymbol{t}_{1:i-1})$

- Incremental initialization is easy to implement in a Gibbs sampler

- Incremental initialization improves token f-score in all models, especially on simple models

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

*but see caveats on next slide!*

MACQUARIE UNIVERSITY

108/117

# Incremental initialization produces low-probability parses

# Why incremental initialization produces low-probability parses

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

# Table label resampling

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

# Table label resampling example

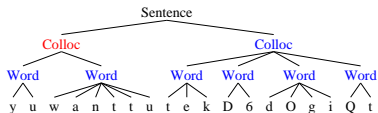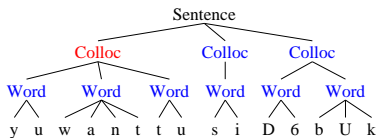Label on table in Chinese Restaurant for colloc

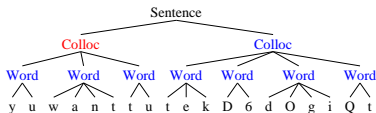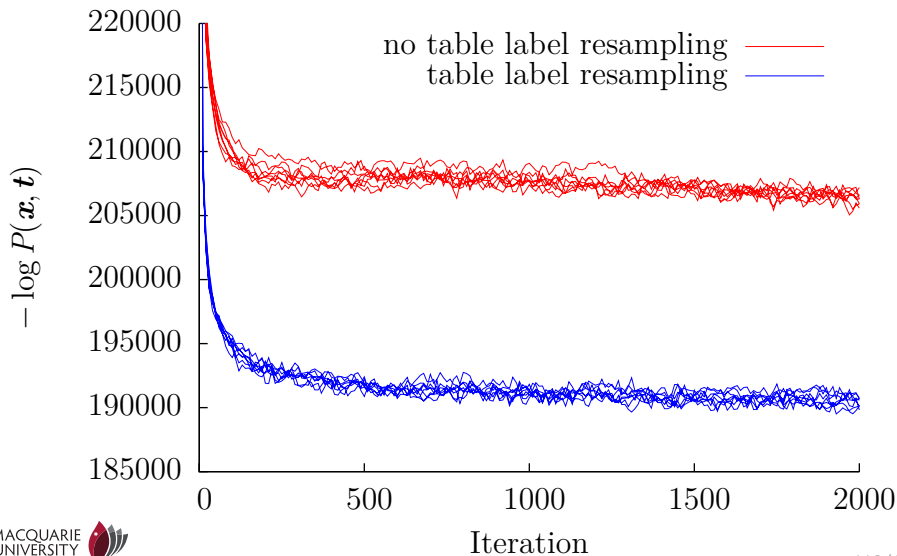

Resulting changes in parse trees

# Table label resampling produces much higher-probability parses

# Summary: learning adaptor grammars

- Unbounded number of possible cached subtrees $\Rightarrow$ Expectation Maximisation isn't sufficient
- *Gibbs sampler* batch learning algorithm
  - assign every sentence a (random) parse
  - repeatedly cycle through training sentences:
    - withdraw parse (decrement counts) for sentence
    - sample parse for current sentence and update counts
    - Metropolis-Hastings correction

# Outline

MACQUARIE
UNIVERSITY

# Conclusions and future work

- Adaptor Grammars can express a variety of useful HDP models
  - generic AG inference code makes it easy to explore models
- AGs have a variety of applications
  - unsupervised acquisition of morphology
  - unsupervised word segmentation
  - learning word to referent mappings
  - learning collocations in topic models
- Future work:
  - extend expressive power of AGs (e.g., feature-passing)
  - richer data (e.g., more non-linguistic context)
  - more realistic data (e.g., phonological variation)

MACQUARIE
UNIVERSITY

# Interested in **statistical models**, **machine learning** and **computational linguistics**?

**Macquarie University** is recruiting
**PhD students** and **post-docs**!

Contact **Mark.Johnson@mq.edu.au** for more information.