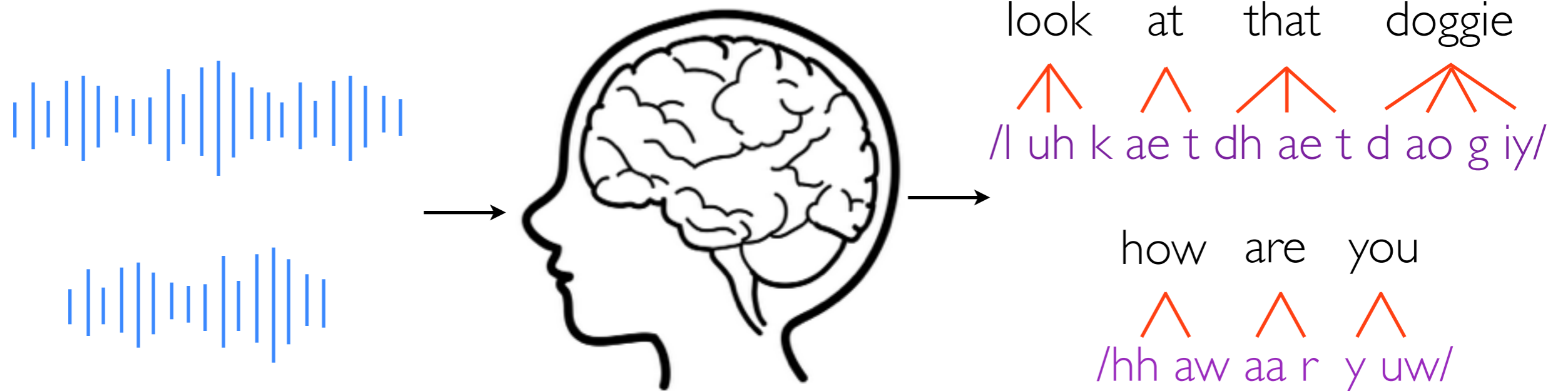# Discovering Linguistic Structures from Speech: Models and Applications

Jackie Lee

Spoken Language Systems Group, CSAIL, MIT
Machine Learning Engineer @ Kite

# Problem Overview
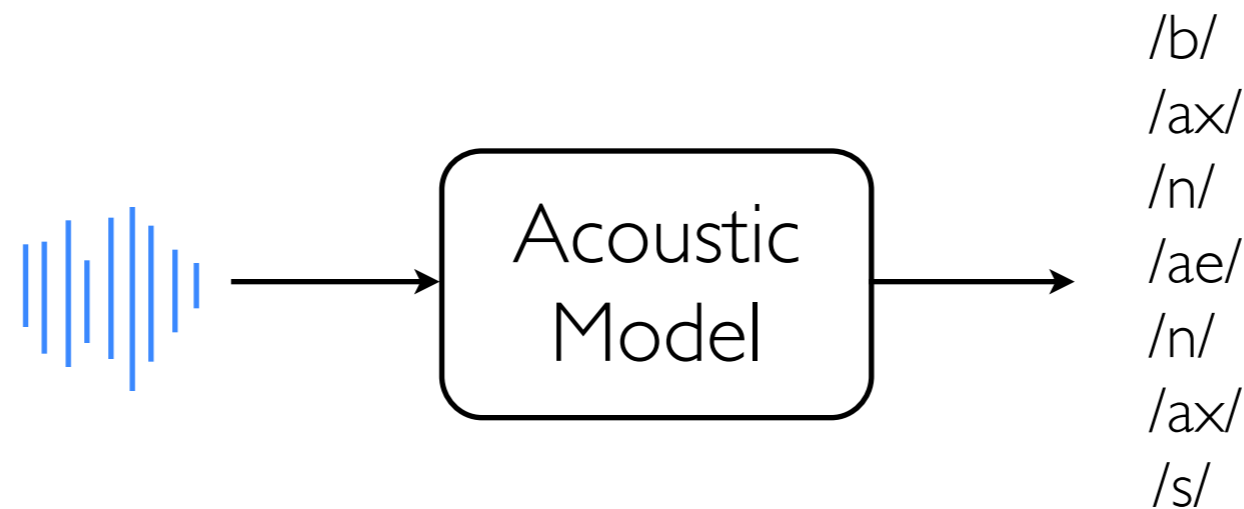
- A task that humans can perform naturally



- Goal

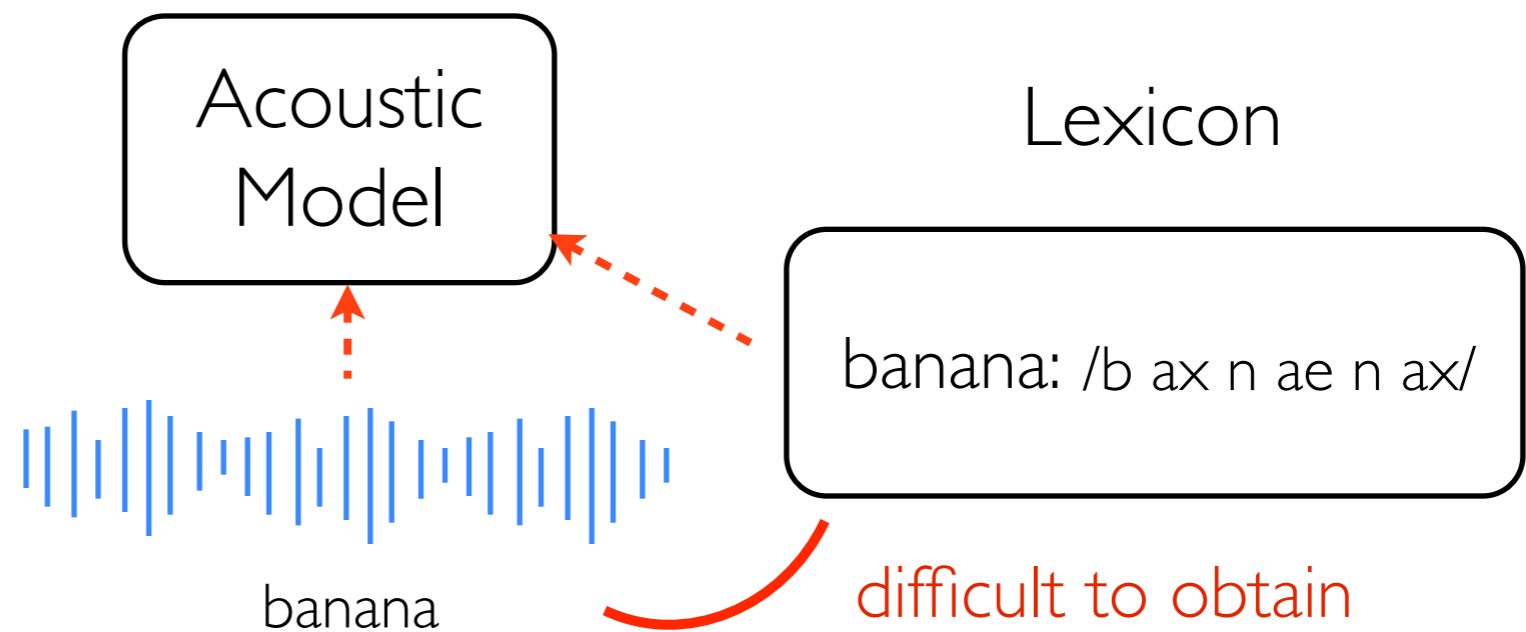    – Develop computational models for discovering linguistic structures from speech

# Potential Applications of Discovered Structures

- Unsupervised training of speech recognizers

- Take acoustic model as an example



/b/
/ax/
/n/
/ae/
/n/
/ax/
/s/

# Potential Applications of Discovered Structures

- Unsupervised training of speech recognizers

- Take acoustic model as an example

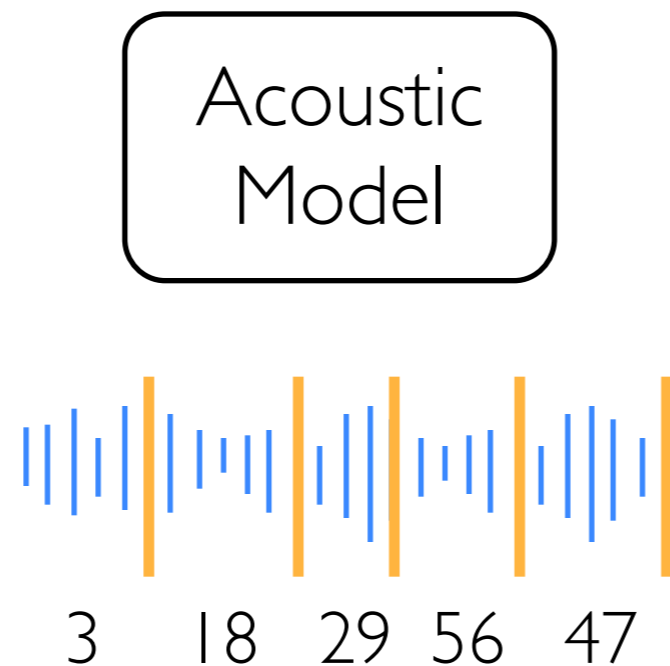  – Training requires word transcriptions with a pronunciation lexicon

# Potential Applications of Discovered Structures

- Unsupervised training of speech recognizers

- Take acoustic model as an example

  - Training requires word transcriptions with a pronunciation lexicon

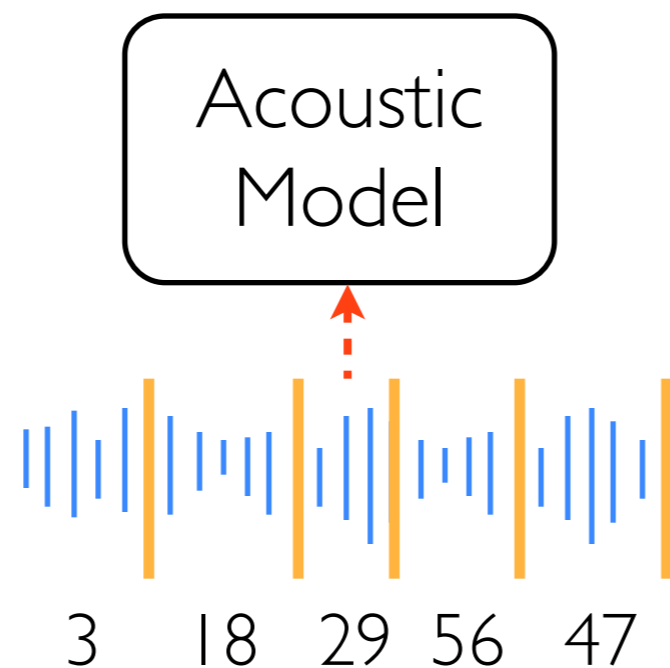- Unsupervised phonetic unit discovery

# Potential Applications of Discovered Structures

- Unsupervised training of speech recognizers

- Take acoustic model as an example

  - Training requires word transcriptions with a pronunciation lexicon

- Unsupervised phonetic unit discovery

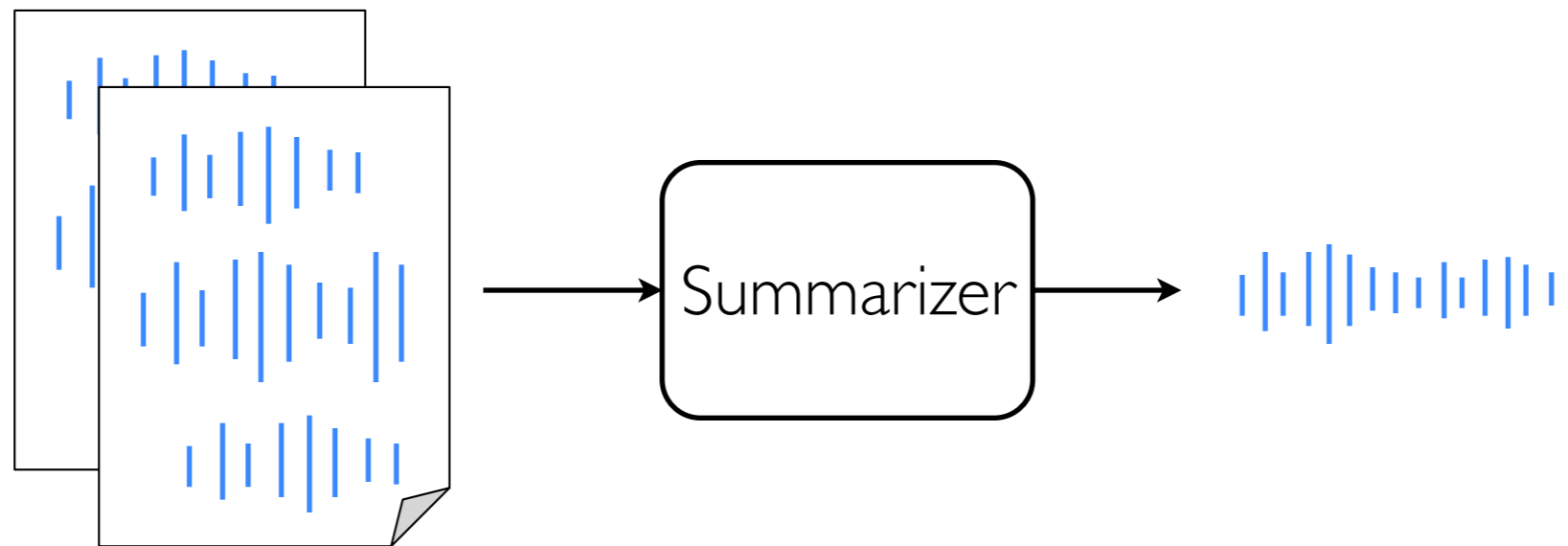  - Allows learning an acoustic model directly from speech data

# Applications of Higher Level Linguistic Structures

- Sub-word units are useful for representing out-of-vocabulary words

# Applications of Higher Level Linguistic Structures

- Sub-word units are useful for representing out-of-vocabulary words

- Unsupervised word discovery

  – Natural language processing on spoken documents without speech recognition



- Connection to the field of Cognitive Science

# Outline

## Discovering phonetic inventory
*[Lee and Glass, ACL 2012]*

/b/  /ax/  /n/  /ae/  /n/  /ax/

Part I of the talk

## Discovering hierarchical linguistic structures
*[Lee, O'Donnell, and Glass, TACL 2015]*

Word      banana

Syllable

Phone   /b/  /ax/  /n/  /ae/  /n/  /ax/

Part II of the talk

# Part I: Discovering Phonetic Units from Speech

**Discovering phonetic inventory**
*[Lee and Glass, ACL 2012]*

/b/  /ax/  /n/  /ae/  /n/  /ax/
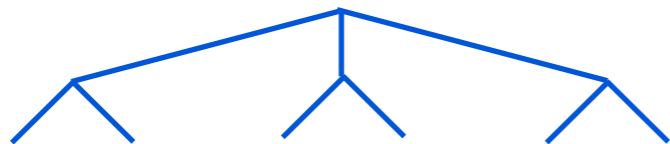
Part I of the talk
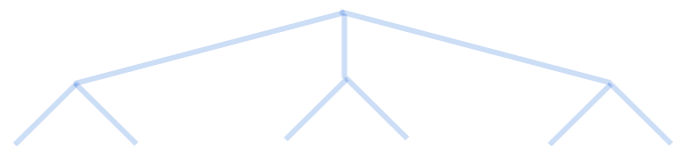
**Discovering hierarchical linguistic structures**
*[Lee, O'Donnell, and Glass, TACL 2015*

Word                    banana

Syllable

Phone   /b/  /ax/  /n/  /ae/  /n/  /ax/
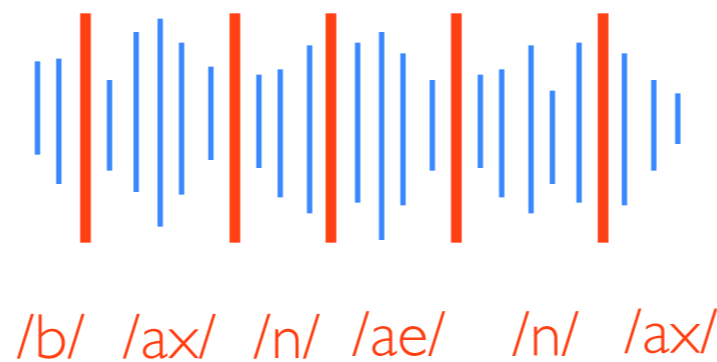
Part II of the talk

# Problem Overview

- Find the phone units embedded in the observed speech data

# Problem Overview

- Find the phone units embedded in the observed speech data

- Latent variables

/b/ /ax/ /n/ /ae/ /n/ /ax/

/b/, /k/, /d/, /ae/, /ix/, /iy/, /e/, /

- Phone boundaries
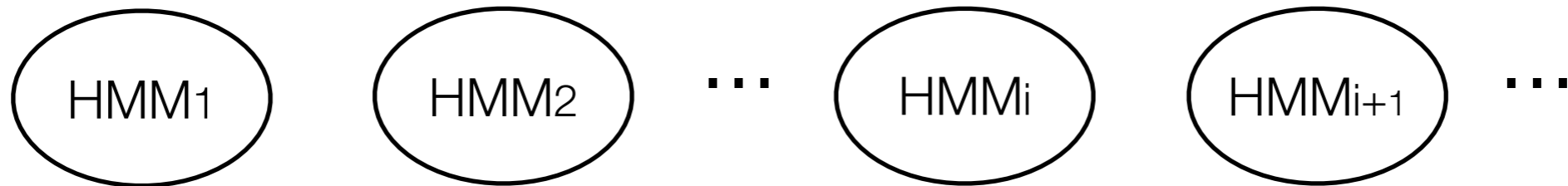
- Phone labels

- Phone inventory

# Related Work

- Unsupervised acoustic unit discovery and modeling

  - Towards unsupervised training of speaker independent acoustic models [*Jansen and Church, INTERSPEECH 2011*]

  - Unsupervised hidden Markov modeling of spoken queries for spoken term detection without speech recognition [*Chan et al., INTERSPEECH 2011*]

  - Keyword spotting of arbitrary words using minimal speech resources [*Garcia and Gish, ICASSP 2006*]

  - Toward ALISP: A proposal for automatic language independent speech processing [*Chollet et al., Computational Models of Speech Pattern Processing 1999*]

  - A segment model based approach to speech recognition [*Lee et al., ICASSP 1988*]

# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated
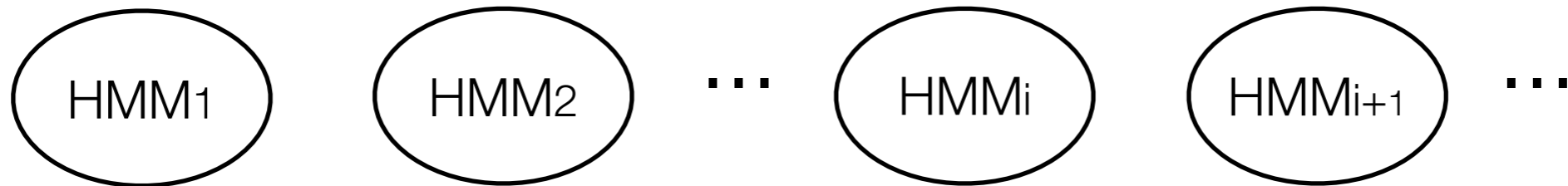
# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated
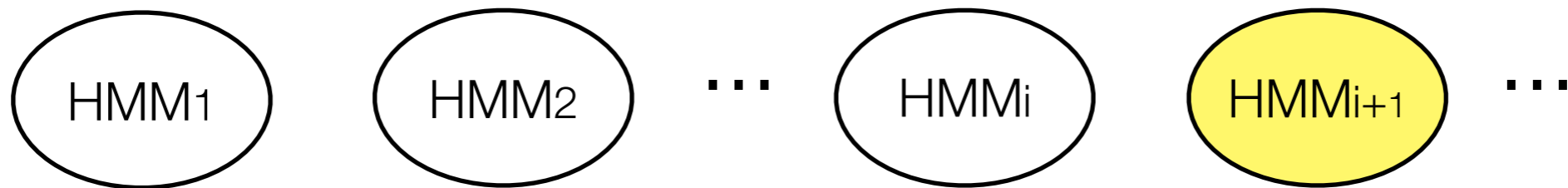
# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated
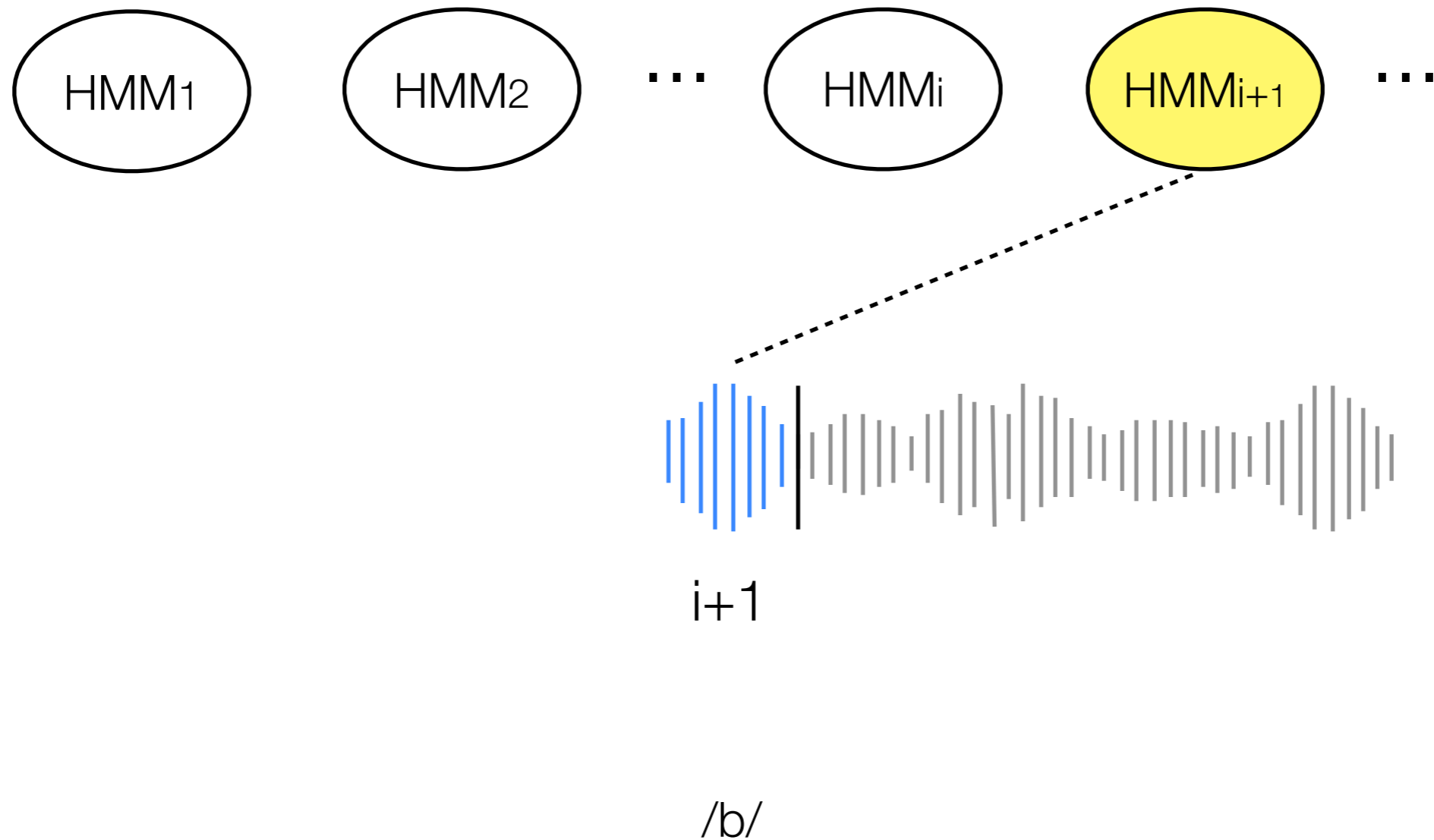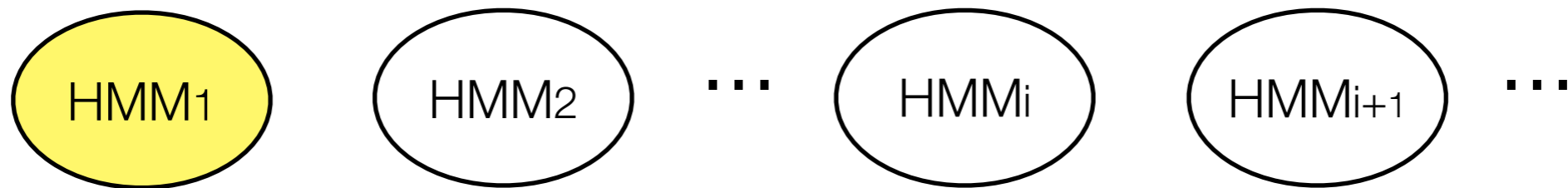
# Generative Story

- A simple explanation of how a spoken utterance is generated

# Generative Story

- A simple explanation of how a spoken utterance is generated
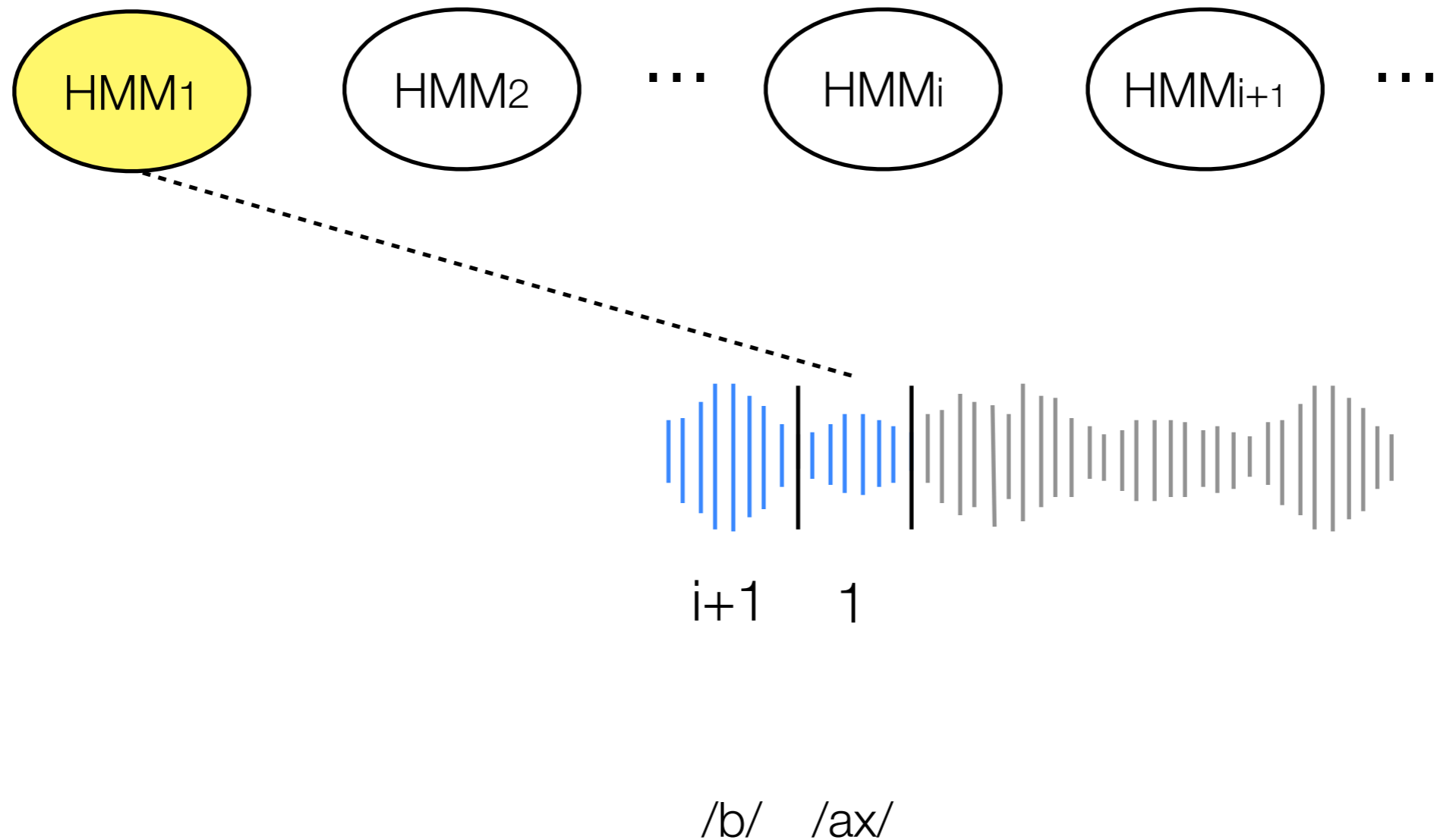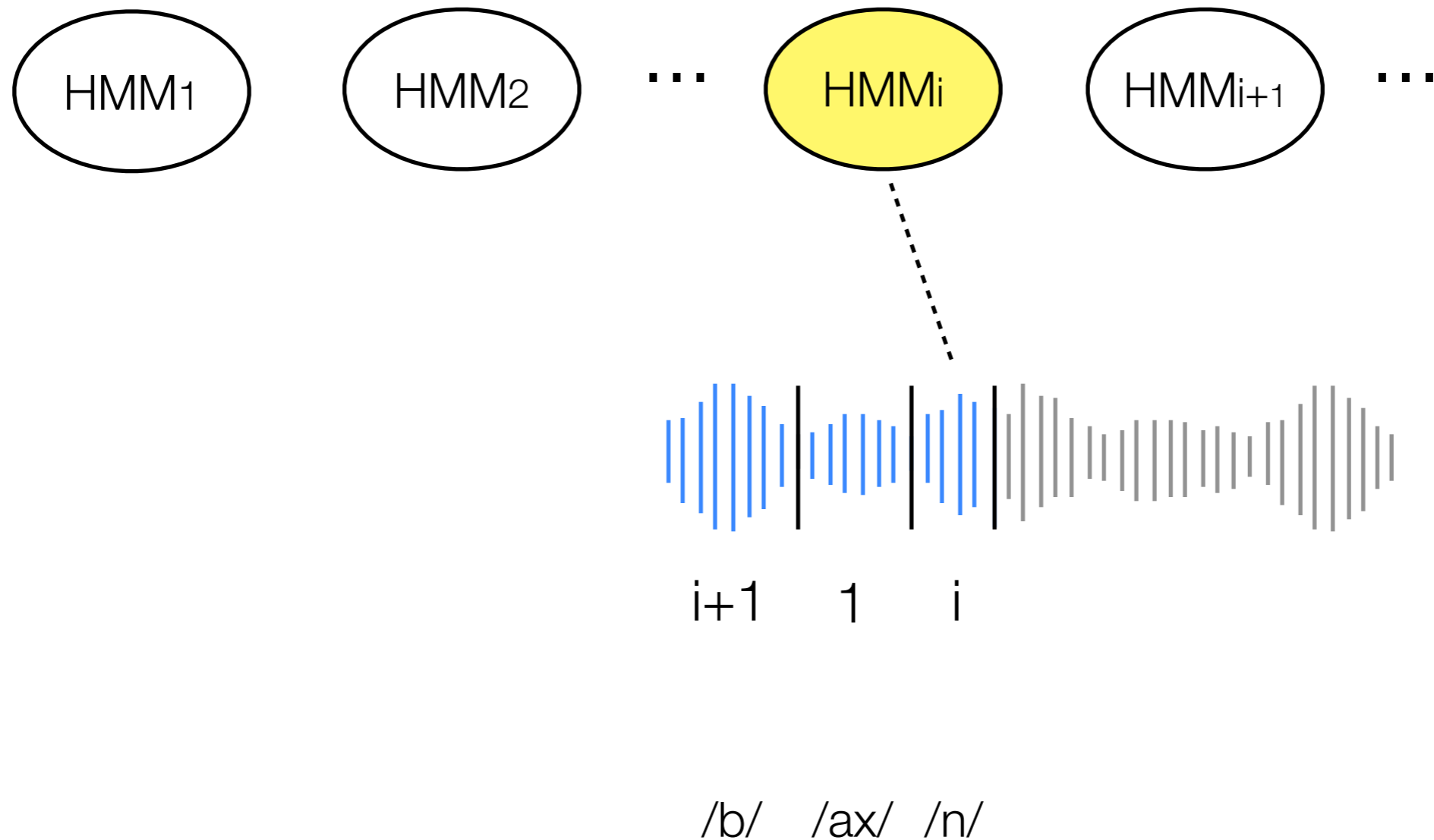
# Generative Story

- A simple explanation of how a spoken utterance is generated

$HMM_1$   $HMM_2$   $\cdots$   $HMM_i$   $HMM_{i+1}$   $\cdots$

- Main latent variables

  – Phone boundaries ($b$)

$b_1$ ... $b_9$ ... $b_{16}$ ... $b_{28}$ ... $b_{37}$ ...

i+1   1   i   2   i   1

/b/   /ax/   /n/   /ae/   /n/   /ax/

# Generative Story

- A simple explanation of how a spoken utterance is generated



- Main latent variables

  - Phone boundaries ($b$)

  - Phone labels ($c$)

# Generative Story

- A simple explanation of how a spoken utterance is generated



- Main latent variables

  - Phone boundaries ($b$)

  - Phone labels ($c$)

  - HMM parameters ($\theta$)

# Generative Story

- A simple explanation of how a spoken utterance is generated



- Main latent variables

  - Phone boundaries (*b*)

  - Phone labels (*c*)

  - HMM parameters (θ)

  - # of HMMs (phones)

  Dirichlet Process

# Inference Procedure

Initialize boundary variables ($b_t$) randomly

$\downarrow$

Sample $c_i$ for each segment

$\downarrow$

Sample HMM parameters ($\theta_i$)

$\downarrow$

Sample for each $b_t$

# Inference Procedure

# Inference Procedure



Initialize boundary variables ($b_t$) randomly

Sample $c_i$ for each segment

Sample HMM parameters ($\theta_i$)

Sample for each $b_t$

# Inference Procedure

Initialize boundary variables ($b_t$) randomly

Sample $c_i$ for each segment

Sample HMM parameters ($\theta_i$)

Sample for each $b_t$

# Inference Procedure

Initialize boundary variables ($b_t$) randomly

Sample $c_i$ for each segment

Sample HMM parameters ($\theta_i$)

Sample for each $b_t$

# Inference Procedure



Initialize boundary variables ($b_t$) randomly

Sample $c_i$ for each segment

Sample HMM parameters ($\theta_i$)

Sample for each $b_t$

Gibbs sampling

- **Iterate n times**
  - n = 20,000 in our experiments

# Inference Procedure



- **Iterate n times**
  - n = 20,000 in our experiments

# DP as a Prior for Phone Labels ($c$)

- A Chinese Restaurant Process (CRP) representation

  - Each table is a phonetic unit

  - Each speech segment is a customer $s_i = [x_t, x_{t+1}, ... x_{t+L_i}]$

# DP as a Prior for Phone Labels ($c$)

- A Chinese Restaurant Process (CRP) representation

  – Each table is a phonetic unit

  – Each speech segment is a customer $s_i = [x_t, x_{t+1}, \ldots x_{t+L_i}]$

$$\boxed{S_1}$$

$$\bigcirc \theta_1$$

$$c_1 = 1$$

# DP as a Prior for Phone Labels (*c*)

- A Chinese Restaurant Process (CRP) representation

  - Each table is a phonetic unit

  - Each speech segment is a customer $s_i = [x_t, x_{t+1}, ... x_{t+Li}]$



$c_1 = 1$      $c_2 = 2$

# DP as a Prior for Phone Labels ($c$)

- A Chinese Restaurant Process (CRP) representation

  - Each table is a phonetic unit

  - Each speech segment is a customer $s_i = [x_t, x_{t+1}, ... x_{t+L_i}]$

$$S_3 \quad S_1 \qquad\qquad S_2$$

$$\theta_1 \qquad\qquad \theta_2$$

$c_1 = 1 \qquad\qquad c_2 = 2$

$c_3 = 1$

# DP as a Prior for Phone Labels ($c$)

- A Chinese Restaurant Process (CRP) representation

  - Each table is a phonetic unit

  - Each speech segment is a customer $s_i = [x_t, x_{t+1}, \ldots x_{t+L_i}]$



$c_1 = 1$  $\qquad$ $c_2 = 2$  $\qquad$ $c_4 = 3$  $\qquad$ $c_t = K$

$c_3 = 1$  $\qquad$ $c_8 = 2$  $\qquad$ $c_5 = 3$

$c_9 = 1$

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

$$p(c_i = k, 1 \le k \le K \mid \cdots) \propto \frac{n_k}{N + \alpha} \, p(s_i \mid \theta_k)$$

$\underbrace{\qquad}_{\text{posterior probability}}$ $\underbrace{\qquad}_{\text{DP prior}}$ $\underbrace{\qquad}_{\text{likelihood}}$

$n_k$ : number of customers at table k

$N$ : number of costumers seen so far

$\alpha$ : concentration parameter of DP

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

$$p(c_i = k, 1 \leq k \leq K \mid \cdots) \propto \frac{n_k}{N + \alpha} \, p(s_i \mid \theta_k)$$

  – $s_i$ opens a new table $\longrightarrow$ $s_i$ is a new phone

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  – $s_i$ sits at an occupied table ⟶ $s_i$ is not a new phone

$$p(c_i = k, 1 \le k \le K \mid \cdots) \propto \frac{n_k}{N + \alpha} p(s_i \mid \theta_k)$$

  – $s_i$ opens a new table ⟶ $s_i$ is a new phone

$$p(c_i = K + 1 \mid \cdots) \propto \frac{\alpha}{N + \alpha} \int_\theta p(s_i \mid \theta) d\theta$$

# Posterior Distribution for $c_i$



- For a new segment ($s_i$), the posterior probability distribution of $c_i$ :

  - $s_i$ sits at an occupied table $\longrightarrow$ $s_i$ is not a new phone

  $$p(c_i = k, 1 \leq k \leq K \mid \cdots) \propto \frac{n_k}{N + \alpha} p(s_i \mid \theta_k)$$

  - $s_i$ opens a new table $\longrightarrow$ $s_i$ is a new phone

  $$p(c_i = K + 1 \mid \cdots) \propto \frac{\alpha}{N + \alpha} \int_\theta p(s_i \mid \theta) d\theta$$

Generate a sample for $c_i$

# Inference Procedure



- Iterate n times

  - n = 20,000 in our experiments

# Inference Procedure



Gibbs sampling

- Initialize boundary variables ($b_t$) randomly
- Sample $c_i$ for each segment
- Sample HMM parameters ($\theta_i$)
- Sample for each $b_t$

- **Iterate n times**
  - n = 20,000 in our experiments

# Experiments

- Data set

  - TIMIT corpus

  - Multi-speaker, clean read speech, 16kHz sampling rate

# Experiments

- Data set

  - TIMIT corpus

  - Multi-speaker, clean read speech, 16kHz sampling rate

- Qualitative assessment

  - Correlation between induced phone units and English phones

  - Results learned from 3696 utterances

# Experiments

- Data set

  - TIMIT corpus

  - Multi-speaker, clean read speech, 16kHz sampling rate

- Qualitative assessment

  - Correlation between induced phone units and English phones

  - Results learned from 3696 utterances

- Quantitative assessments

  - Phone segmentation

  - (Query-by-example spoken term detection)

# Discovered Phone Units -- 3696 utterances

- **123** phone units discovered from **3696 TIMIT** utterances

  - A fine correlation between discovered phones and English phones

# Discovered Phone Units -- 3696 utterances

- **123 phone units discovered from 3696 TIMIT utterances**

  - A fine correlation between discovered phones and English phones

# Discovered Phone Units -- 3696 utterances

- **123** phone units discovered from **3696 TIMIT** utterances

  - A fine correlation between discovered phones and English phones

# Discovered Phone Units -- 3696 utterances

- **123 phone units discovered from 3696 TIMIT utterances**

  - A fine correlation between discovered phones and English phones

# Discovered Phone Units -- 3696 utterances

- **123 phone units discovered from 3696 TIMIT utterances**

  - A fine correlation between discovered phones and English phones



Context-dependent:
/ae/ + /m/, /n/
/ae/ + stops

# Phone Segmentation

- TIMIT training portion

|  | Recall | Precision | F-score |
|---|---|---|---|
| Dusan et al. (unsupervised) | 75.2 | 66.8 | 70.8 |
| Qiao et al. (semi-supervised) | 77.5 | 76.3 | 76.9 |
| Our model (unsupervised) | 76.2 | 76.4 | 76.3 |

# Part I: Discovering Phonetic Units from Speech

## Discovering phonetic inventory
*[Lee and Glass, ACL 2012]*

/b/  /ax/  /n/  /ae/  /n/  /ax/

- DP mixture models with HMMs

  - Discovered phonetic units are highly correlated with standard phones

  - Achieves phone segmentation performance similar to the semi-supervised baseline

## Discovering hierarchical linguistic structures
*[Lee, O'donnell, and Glass, TACL 2015*

Word                    banana

Syllable

Phone    /b/  /ax/  /n/  /ae/  /n/  /ax/

Part II of the talk

# Part II: Discovering Hierarchical Linguistic Structures

Discovering phonetic inventory
[Lee and Glass, ACL 2012]

/b/ /ax/ /n/ /ae/ /n/ /ax/

- DP mixture models with HMMs

  – Discovered phonetic units are highly correlated with standard phones

  – Achieves phone segmentation performance similar to the semi-supervised baseline

## Discovering hierarchical linguistic structures

*[Lee, O'Donnell, and Glass, TACL 2015]*

Word        banana

Syllable

Phone   /b/ /ax/ /n/ /ae/ /n/ /ax/

**Part II of the talk**

# Problem Overview

- Discover hierarchical linguistic structures from speech

  – Phone-like, syllable-like and word-like units

# Related Work

- **Spoken term discovery**

  - Unsupervised patter discovery in speech [*Park and Glass, IEEE Trans., 2008*]

  - Unsupervised speech processing with applications to query-by-example spoken term detection [*Zhang, Ph.D. Thesis 2013*]

  - Towards spoken term discovery at scale with zero resources [*Jansen et al., INTERSPEECH 2010*]

- **Word segmentation on phone transcripts of spoken utterances**

  - A Bayesian framework for word segmentation: Exploring the effects of context [*Goldwater et al., Cognition 2009*]

  - Bayesian unsupervised word segmentation with nested Pitman-Yor language modeling [*Mochihashi et al., ACL 2009*]

  - Using adaptor grammars to identify synergies in the unsupervised acquisition of linguistic structure [*Johnson, ACL-HLT 2008*]

# Spoken Term Discovery

- Discover speech segments that correspond to words

**Input:**

[*Park and Glass, IEEE Trans., 2008*] [*Zhang, Ph.D. Thesis 2013*] [*Jansen et al., INTERSPEECH 2010*]

# Spoken Term Discovery

- Discover speech segments that correspond to words

university

Input: 

[*Park and Glass, IEEE Trans., 2008*] [*Zhang, Ph.D. Thesis 2013*] [*Jansen et al., INTERSPEECH 2010*]

# Word Segmentation on Phone Transcripts

- Model words as sequences of phones

**Words:**     and          MIT's          open                university                      and

**Input:**  ae  n  d  eh  m  ay  t  iy  z  ow  p  ax  n  y  uw  n  ax  v  er  s  ax  dx  iy  ae  n  d

[*Goldwater et al., ACL 2006*]  [*Brent and Cartwrite, Cognition 1996*]  [*Mochihashi et al., ACL 2009*]

# Word Segmentation on Phone Transcripts

- Model words as sequences of phones

- Modeling more levels of structures improves word segmentation

  - Word → Syllables        Syllable → Phones

**Words:**  and    MIT's    open    university    and

**Input:**  ae n d  eh m ay t iy z  ow p ax n  y uw n ax v er s ax dx iy  ae n d

[*Johnson, ACL-HLT 2008*]  [*Johnson et al., NAACL-HLT, 2009*]

# Word Segmentation on Phone Transcripts

- Model words as sequences of phones

- Modeling more levels of structures improves word segmentation

  – Word → Syllables          Syllable → Phones

**Words:**     and          MIT's        open              university              and

**Syllables:** [ae n d] [eh m] [ay] [t iy z] [ow p] [ax n] [y uw] [n ax] [v er] [s ax] [dx iy] [ae n d]

**Input:**     ae n d eh m ay t iy z ow p ax n y uw n ax v er s ax dx iy ae n d

[*Johnson, ACL-HLT 2008*]  [*Johnson et al., NAACL-HLT, 2009*]

# Word Segmentation on Phone Transcripts

- Model words as sequences of phones

- Modeling more levels of structures improves word segmentation

  – Word → Syllables        Syllable → Phones

- Adaptor grammars is an effective tool for learning rich structures



**Words:**  and    MIT's    open    university    and

**Syllables:** [ae n d] [eh m] [ay] [t iy z] [ow p] [ax n] [y uw] [n ax] [v er] [s ax] [dx iy] [ae n d]

**Input:**  ae n d eh m ay t iy z ow p ax n y uw n ax v er s ax dx iy ae n d

only learns from symbolic input

[*Johnson, ACL-HLT 2008*]  [*Johnson et al., NAACL-HLT, 2009*]

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phone discovery model

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  – To discover rich linguistic structures from speech

- Three components in the model

Adaptor grammars

A nonparametric Bayesian extension of probabilistic context-free grammars (PCFGs)

Noisy-channel model

Phone discovery model

# PCFG Example

An example PCFG for
generating phone sequences

PCFG

| | | |
|---|---|---|
| 0.5 | Sen $\longrightarrow$ | Word Word |
| 0.5 | Sen $\longrightarrow$ | Word |
| 0.7 | Word $\longrightarrow$ | Syl Syl |
| 0.3 | Word $\longrightarrow$ | Syl |
| 1.0 | Syl $\longrightarrow$ | Phn Phn |
| 0.1 | Phn $\longrightarrow$ | /ax/ |
| 0.05 | Phn $\longrightarrow$ | /n/ |
| 0.1 | Phn $\longrightarrow$ | /ow/ |
| 0.1 | Phn $\longrightarrow$ | /p/ |
| | ... | |

# PCFG Generative Process

Sen

**PCFG**

| 0.5 | Sen | $\longrightarrow$ | Word Word |
| 0.5 | Sen | $\longrightarrow$ | Word |
| 0.7 | Word | $\longrightarrow$ | Syl Syl |
| 0.3 | Word | $\longrightarrow$ | Syl |
| 1.0 | Syl | $\longrightarrow$ | Phn Phn |
| 0.1 | Phn | $\longrightarrow$ | /ax/ |
| 0.05 | Phn | $\longrightarrow$ | /n/ |
| 0.1 | Phn | $\longrightarrow$ | /ow/ |
| 0.1 | Phn | $\longrightarrow$ | /p/ |
| | | | ... |

# PCFG Generative Process

Sen
|
Word

### PCFG

| | | |
|---|---|---|
| 0.5 | Sen | ⟶ Word Word |
| 0.5 | Sen | ⟶ Word |
| 0.7 | Word | ⟶ Syl Syl |
| 0.3 | Word | ⟶ Syl |
| 1.0 | Syl | ⟶ Phn Phn |
| 0.1 | Phn | ⟶ /ax/ |
| 0.05 | Phn | ⟶ /n/ |
| 0.1 | Phn | ⟶ /ow/ |
| 0.1 | Phn | ⟶ /p/ |
| | | ... |

# PCFG Generative Process



Sen

Word

Syl    Syl

PCFG

| 0.5 | Sen | $\longrightarrow$ | Word Word |
| 0.5 | Sen | $\longrightarrow$ | Word |
| 0.7 | Word | $\longrightarrow$ | Syl Syl |
| 0.3 | Word | $\longrightarrow$ | Syl |
| 1.0 | Syl | $\longrightarrow$ | Phn Phn |
| 0.1 | Phn | $\longrightarrow$ | /ax/ |
| 0.05 | Phn | $\longrightarrow$ | /n/ |
| 0.1 | Phn | $\longrightarrow$ | /ow/ |
| 0.1 | Phn | $\longrightarrow$ | /p/ |

...

# PCFG Generative Process

# PCFG Generative Process



PCFG

| | | |
|---|---|---|
| 0.5 | Sen | ⟶ Word Word |
| 0.5 | Sen | ⟶ Word |
| 0.7 | Word | ⟶ Syl Syl |
| 0.3 | Word | ⟶ Syl |
| 1.0 | Syl | ⟶ Phn Phn |
| 0.1 | Phn | ⟶ /ax/ |
| 0.05 | Phn | ⟶ /n/ |
| 0.1 | Phn | ⟶ /ow/ |
| 0.1 | Phn | ⟶ /p/ |
| | ... | |

# PCFG Generative Process

# PCFG Generative Process



Sen

Word

Syl     Syl

Phn Phn Phn Phn

/ow/   /p/   /ax/   /n/

**PCFG**

| 0.5 | Sen | ⟶ | Word Word |
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| ... |

# Adaptor Grammars

- A PCFG +

PCFG

| | | |
|---|---|---|
| 0.5 | Sen | ⟶ Word Word |
| 0.5 | Sen | ⟶ Word |
| 0.7 | Word | ⟶ Syl Syl |
| 0.3 | Word | ⟶ Syl |
| 1.0 | Syl | ⟶ Phn Phn |
| 0.1 | Phn | ⟶ /ax/ |
| 0.05 | Phn | ⟶ /n/ |
| 0.1 | Phn | ⟶ /ow/ |
| 0.1 | Phn | ⟶ /p/ |
| | ... | |

# Adaptor Grammars

- A PCFG + cached subtrees for adapted nonterminals

PCFG

| 0.5 | Sen | → | Word Word |
|---|---|---|---|
| 0.5 | Sen | → | Word |
| 0.7 | Word | → | Syl Syl |
| 0.3 | Word | → | Syl |
| 1.0 | Syl | → | Phn Phn |
| 0.1 | Phn | → | /ax/ |
| 0.05 | Phn | → | /n/ |
| 0.1 | Phn | → | /ow/ |
| 0.1 | Phn | → | /p/ |

...

Cached subtrees

Syl:

# Adaptor Grammars

- A PCFG + cached subtrees for adapted nonterminals

Key idea:
Adaptor grammars memorize
reusable structures

PCFG

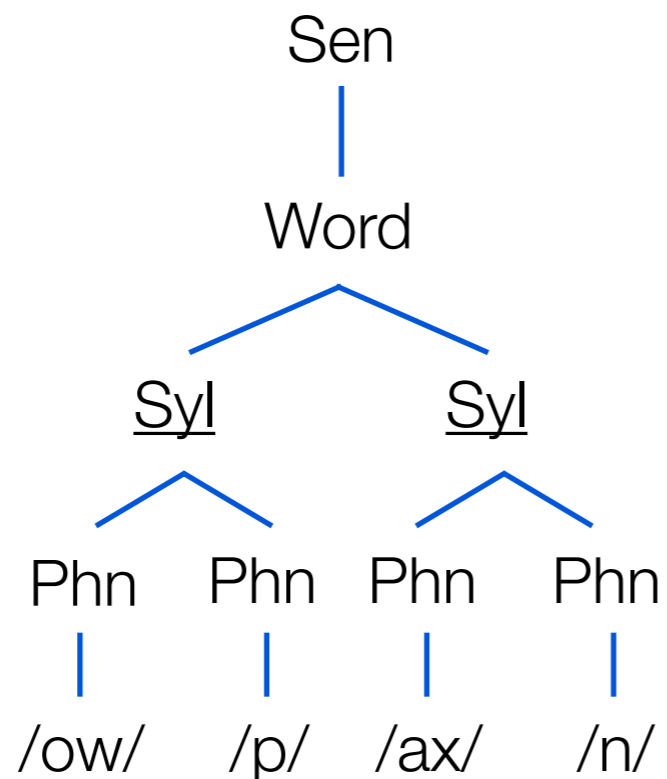| 0.5 | Sen | ⟶ | Word Word |
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | | ... |

Cached
subtrees

Syl:

# Adaptor Grammars Generative Process

- Assume a current parse



Cached subtrees

Syl:

PCFG

| 0.5 | Sen | ⟶ | Word Word |
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | <u>Syl</u> <u>Syl</u> |
| 0.3 | Word | ⟶ | <u>Syl</u> |
| 1.0 | <u>Syl</u> | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | ... | |

# Adaptor Grammars Generative Process

- Cache subtrees for adapted nonterminals

# Adaptor Grammars Generative Process

- Cache subtrees for adapted nonterminals



PCFG

| | | |
|---|---|---|
| 0.5 | Sen | ⟶ Word Word |
| 0.5 | Sen | ⟶ Word |
| 0.7 | Word | ⟶ Syl Syl |
| 0.3 | Word | ⟶ Syl |
| 1.0 | Syl | ⟶ Phn Phn |
| 0.1 | Phn | ⟶ /ax/ |
| 0.05 | Phn | ⟶ /n/ |
| 0.1 | Phn | ⟶ /ow/ |
| 0.1 | Phn | ⟶ /p/ |
| | ... | |

# Adaptor Grammars Generative Process

- Generate a new parse

Sen

PCFG

| 0.5 | Sen | ⟶ | Word Word |
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | | ... |

Cached subtrees

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- Expand regular nonterminals using PCFG rules

Sen

Word

PCFG

| 0.5 | Sen | → | Word Word |
| 0.5 | Sen | → | Word |
| 0.7 | Word | → | <u>Syl</u> <u>Syl</u> |
| 0.3 | Word | → | <u>Syl</u> |
| 1.0 | <u>Syl</u> | → | Phn Phn |
| 0.1 | Phn | → | /ax/ |
| 0.05 | Phn | → | /n/ |
| 0.1 | Phn | → | /ow/ |
| 0.1 | Phn | → | /p/ |

...

Cached subtrees

<u>Syl</u>:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- Expand regular nonterminals using PCFG rules

Sen

(Word)

Syl

PCFG

| 0.5 | Sen | → | Word Word |
| 0.5 | Sen | → | Word |
| 0.7 | Word | → | Syl Syl |
| 0.3 | Word | → | Syl |
| 1.0 | Syl | → | Phn Phn |
| 0.1 | Phn | → | /ax/ |
| 0.05 | Phn | → | /n/ |
| 0.1 | Phn | → | /ow/ |
| 0.1 | Phn | → | /p/ |
| | | | ... |

Cached subtrees

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

Sen
|
Word
|
(Syl)

### PCFG

| | | |
|---|---|---|
| 0.5 | Sen ⟶ | Word Word |
| 0.5 | Sen ⟶ | Word |
| 0.7 | Word ⟶ | Syl Syl |
| 0.3 | Word ⟶ | Syl |
| 1.0 | Syl ⟶ | Phn Phn |
| 0.1 | Phn ⟶ | /ax/ |
| 0.05 | Phn ⟶ | /n/ |
| 0.1 | Phn ⟶ | /ow/ |
| 0.1 | Phn ⟶ | /p/ |
| | | ... |

Cached subtrees

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

$$p\left(\begin{array}{c}\text{Syl}\\ \text{/ow/ /p/}\end{array}\right) = \frac{1}{2 + \alpha_{syl}}$$

$$p\left(\begin{array}{c}\text{Syl}\\ \text{/ax/ /n/}\end{array}\right) = \frac{1}{2 + \alpha_{syl}}$$

Sen

|

Word

|

(Syl)

**PCFG**

| | | |
|---|---|---|
| 0.5 | Sen $\longrightarrow$ | Word Word |
| 0.5 | Sen $\longrightarrow$ | Word |
| 0.7 | Word $\longrightarrow$ | Syl Syl |
| 0.3 | Word $\longrightarrow$ | Syl |
| 1.0 | Syl $\longrightarrow$ | Phn Phn |
| 0.1 | Phn $\longrightarrow$ | /ax/ |
| 0.05 | Phn $\longrightarrow$ | /n/ |
| 0.1 | Phn $\longrightarrow$ | /ow/ |
| 0.1 | Phn $\longrightarrow$ | /p/ |
| | ... | |

**Cached subtrees**

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

$$p\left( \overset{\text{Syl}}{\underset{/ow/\ /p/}{\wedge}} \right) = \frac{1}{2 + \alpha_{syl}}$$

$$p\left( \overset{\text{Syl}}{\underset{/ax/\ /n/}{\wedge}} \right) = \frac{1}{2 + \alpha_{syl}}$$

  - Create and store a new subtree

Sen

Word

Syl

**PCFG**

| | | |
|---|---|---|
| 0.5 | Sen $\longrightarrow$ | Word Word |
| 0.5 | Sen $\longrightarrow$ | Word |
| 0.7 | Word $\longrightarrow$ | Syl Syl |
| 0.3 | Word $\longrightarrow$ | Syl |
| 1.0 | Syl $\longrightarrow$ | Phn Phn |
| 0.1 | Phn $\longrightarrow$ | /ax/ |
| 0.05 | Phn $\longrightarrow$ | /n/ |
| 0.1 | Phn $\longrightarrow$ | /ow/ |
| 0.1 | Phn $\longrightarrow$ | /p/ |
| | ... | |

Cached subtrees

Syl:

Syl
/ow//p/

Syl
/ax//n/

Syl

# Adaptor Grammars Generative Process

- Expand adapted nonterminals

  – Reuse a cached subtree

$$p(\ \overset{\text{Syl}}{\underset{/\text{ow}/\ /\text{p}/}{\wedge}}\ ) = \frac{1}{2 + \alpha_{syl}}$$

$$p(\ \overset{\text{Syl}}{\underset{/\text{ax}/\ /\text{n}/}{\wedge}}\ ) = \frac{1}{2 + \alpha_{syl}}$$

  – Create and store a
    new subtree

$$p(\ \overset{\text{Syl}}{\wedge}\ ) = \frac{\alpha_{syl}}{2 + \alpha_{syl}}$$

Sen

|
Word
|
(Syl)

**PCFG**

| | | |
|---|---|---|
| 0.5 | Sen | $\longrightarrow$ Word Word |
| 0.5 | Sen | $\longrightarrow$ Word |
| 0.7 | Word | $\longrightarrow$ Syl Syl |
| 0.3 | Word | $\longrightarrow$ Syl |
| 1.0 | Syl | $\longrightarrow$ Phn Phn |
| 0.1 | Phn | $\longrightarrow$ /ax/ |
| 0.05 | Phn | $\longrightarrow$ /n/ |
| 0.1 | Phn | $\longrightarrow$ /ow/ |
| 0.1 | Phn | $\longrightarrow$ /p/ |
| | | ... |

Cached subtrees

Syl:

Syl
/ow//p/   •

Syl
/ax//n/   •

Syl
∧

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

  $$p\left( \overset{\text{Syl}}{\underset{/ow/\ /p/}{\wedge}} \right) = \frac{1}{2 + \alpha_{syl}}$$

  $$p\left( \overset{\text{Syl}}{\underset{/ax/\ /n/}{\wedge}} \right) = \frac{1}{2 + \alpha_{syl}}$$

  - Create and store a new subtree

  $$p\left( \overset{\text{Syl}}{\wedge} \right) = \frac{\alpha_{syl}}{2 + \alpha_{syl}}$$

**Sen**
|
**Word**
|
(**Syl**)

Cached subtrees

Syl:  Syl /ow//p/  •   Syl /ax//n/  •   Syl

### PCFG

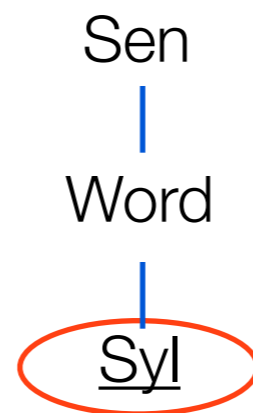| 0.5   | Sen  | → | Word Word |
|-------|------|---|-----------|
| 0.5   | Sen  | → | Word      |
| 0.7   | Word | → | Syl Syl   |
| 0.3   | Word | → | Syl       |
| 1.0   | Syl  | → | Phn Phn   |
| 0.1   | Phn  | → | /ax/      |
| 0.05  | Phn  | → | /n/       |
| 0.1   | Phn  | → | /ow/      |
| 0.1   | Phn  | → | /p/       |
|       |      |   | ...       |

# Adaptor Grammars Generative Process

- Expand adapted nonterminals

  - Reuse a cached subtree

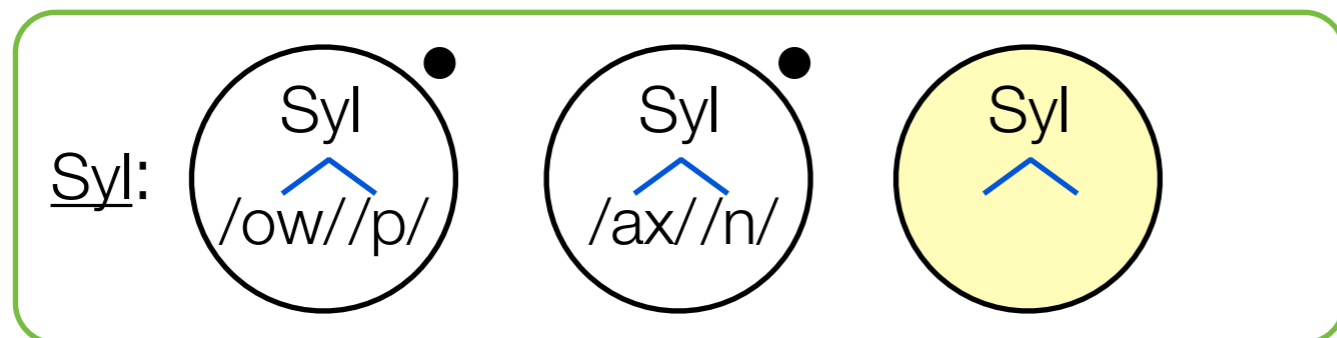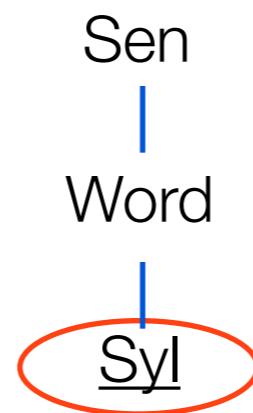  $$p\left(\begin{array}{c}\text{Syl}\\ \text{/ow/ /p/}\end{array}\right) = \frac{1}{2 + \alpha_{\text{syl}}}$$
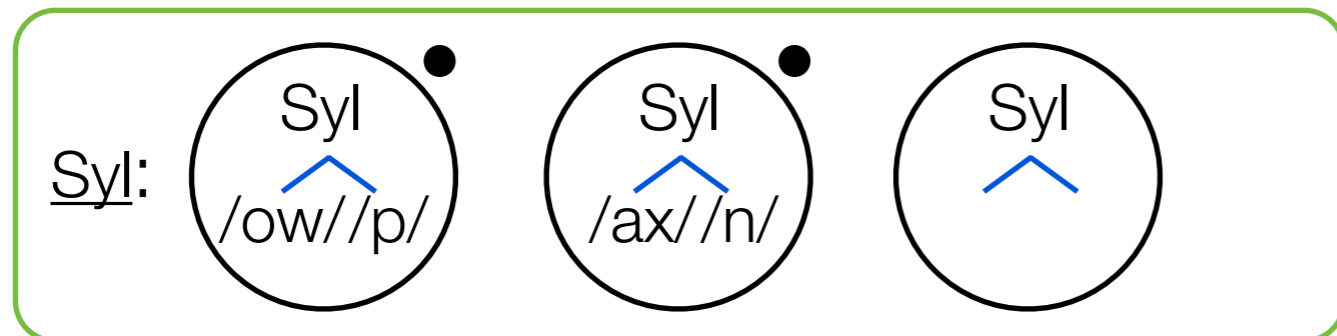
  $$p\left(\begin{array}{c}\text{Syl}\\ \text{/ax/ /n/}\end{array}\right) = \frac{1}{2 + \alpha_{\text{syl}}}$$

  - Create and store a new subtree

  $$p\left(\begin{array}{c}\text{Syl}\\ \wedge\end{array}\right) = \frac{\alpha_{\text{syl}}}{2 + \alpha_{\text{syl}}}$$

Sen
|
Word
|
Syl

Cached subtrees

Syl: Syl /ow//p/ • | Syl /ax//n/ • | Syl

## PCFG

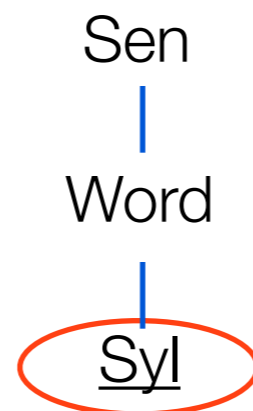| 0.5 | Sen | ⟶ | Word Word |
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | | ... |

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

$$p\left(\begin{array}{c}\text{Syl}\\ \text{/ow/ /p/}\end{array}\right) = \frac{1}{2 + \alpha_{syl}}$$

$$p\left(\begin{array}{c}\text{Syl}\\ \text{/ax/ /n/}\end{array}\right) = \frac{1}{2 + \alpha_{syl}}$$

  - Create and store a new subtree

$$p\left(\begin{array}{c}\text{Syl}\\ \wedge\end{array}\right) = \frac{\alpha_{syl}}{2 + \alpha_{syl}}$$

Sen

|

Word

|

(Syl)

Cached subtrees

Syl:

| Syl /ow//p/ ● | Syl /ax//n/ ● | Syl |

### PCFG

| 0.5 | Sen | ⟶ | Word Word |
|---|---|---|---|
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | ... | |

# Adaptor Grammars Generative Process

- **Expand adapted nonterminals**

  - Reuse a cached subtree

$$p\left(\;\overset{\text{Syl}}{\underset{\text{/ow/ /p/}}{\wedge}}\;\right) = \frac{1}{2 + \alpha_{\text{syl}}}$$

$$p\left(\;\overset{\text{Syl}}{\underset{\text{/ax/ /n/}}{\wedge}}\;\right) = \frac{1}{2 + \alpha_{\text{syl}}}$$

  - Create and store a new subtree

$$p\left(\;\overset{\text{Syl}}{\wedge}\;\right) = \frac{\alpha_{\text{syl}}}{2 + \alpha_{\text{syl}}}$$

Cached subtrees

Syl:

Sen
|
Word
|
(Syl)
/\
/ax/ /n/

PCFG

| 0.5 | Sen | → | Word Word |
| 0.5 | Sen | → | Word |
| 0.7 | Word | → | Syl Syl |
| 0.3 | Word | → | Syl |
| 1.0 | Syl | → | Phn Phn |
| 0.1 | Phn | → | /ax/ |
| 0.05 | Phn | → | /n/ |
| 0.1 | Phn | → | /ow/ |
| 0.1 | Phn | → | /p/ |
| | | | ... |

# Adaptor Grammars Generative Process

- Expand adapted nonterminals

Sen

Word

Syl

/ax/ /n/

PCFG

| | | |
|---|---|---|
| 0.5 | Sen | → Word Word |
| 0.5 | Sen | → Word |
| 0.7 | Word | → Syl Syl |
| 0.3 | Word | → Syl |
| 1.0 | Syl | → Phn Phn |
| 0.1 | Phn | → /ax/ |
| 0.05 | Phn | → /n/ |
| 0.1 | Phn | → /ow/ |
| 0.1 | Phn | → /p/ |
| | | ... |

Cached subtrees

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

# Adaptor Grammars Generative Process

- Cache subtrees for adapted nonterminals

Sen

Word

Syl

/ax/ /n/

Memorize and reuse structures for expanding Syl

Cached subtrees

Syl:

Syl
/ow/ /p/

Syl
/ax/ /n/

PCFG

| 0.5 | Sen | ⟶ | Word Word |
|-----|------|------|------|
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | | ... |

# For Our Problem

- The phone inventory is unknown
  - Terminal symbols should be discovered phonetic unit ids

Sen

Word

Syl

/ax/ /n/

**PCFG**

| 0.5 | Sen | ⟶ | Word Word |
|---|---|---|---|
| 0.5 | Sen | ⟶ | Word |
| 0.7 | Word | ⟶ | Syl Syl |
| 0.3 | Word | ⟶ | Syl |
| 1.0 | Syl | ⟶ | Phn Phn |
| 0.1 | Phn | ⟶ | /ax/ |
| 0.05 | Phn | ⟶ | /n/ |
| 0.1 | Phn | ⟶ | /ow/ |
| 0.1 | Phn | ⟶ | /p/ |
| | | | ... |

Cached subtrees

Syl:

Syl
/ow//p/

Syl
/ax//n/

# For Our Problem

- The phone inventory is unknown
  - Terminal symbols should be discovered phonetic unit ids

Sen

Word

Syl

26  58

PCFG

| 0.5 | Sen | $\longrightarrow$ | Word Word |
| 0.5 | Sen | $\longrightarrow$ | Word |
| 0.7 | Word | $\longrightarrow$ | Syl Syl |
| 0.3 | Word | $\longrightarrow$ | Syl |
| 1.0 | Syl | $\longrightarrow$ | Phn Phn |
| 0.1 | Phn | $\longrightarrow$ | 26 |
| 0.05 | Phn | $\longrightarrow$ | 58 |
| 0.1 | Phn | $\longrightarrow$ | 3 |
| 0.1 | Phn | $\longrightarrow$ | 16 |
|  |  |  | ... |

Cached subtrees

Syl:

Syl
3  16

Syl
26  58

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  – To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phone discovery model

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phone discovery model    First part of the talk

# Recall

- A standard phone may map to multiple discovered units

- Various phone sequences for a word type

# Recall

- A standard phone may map to multiple discovered units

- Various phone sequences for a word type

/k/ /ae/ /t/      /k/ /ae/ /t/

49 (58) 32      49 (26) 32

# Recall

- A standard phone may map to multiple discovered units

- Various phone sequences for a word type

- These variations must be collapsed for lexicon learning

/k/ /ae/ /t/      /k/ /ae/ /t/

49 (58) 32      49 (26) 32

Collapse the variations by
using a noisy-channel model

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

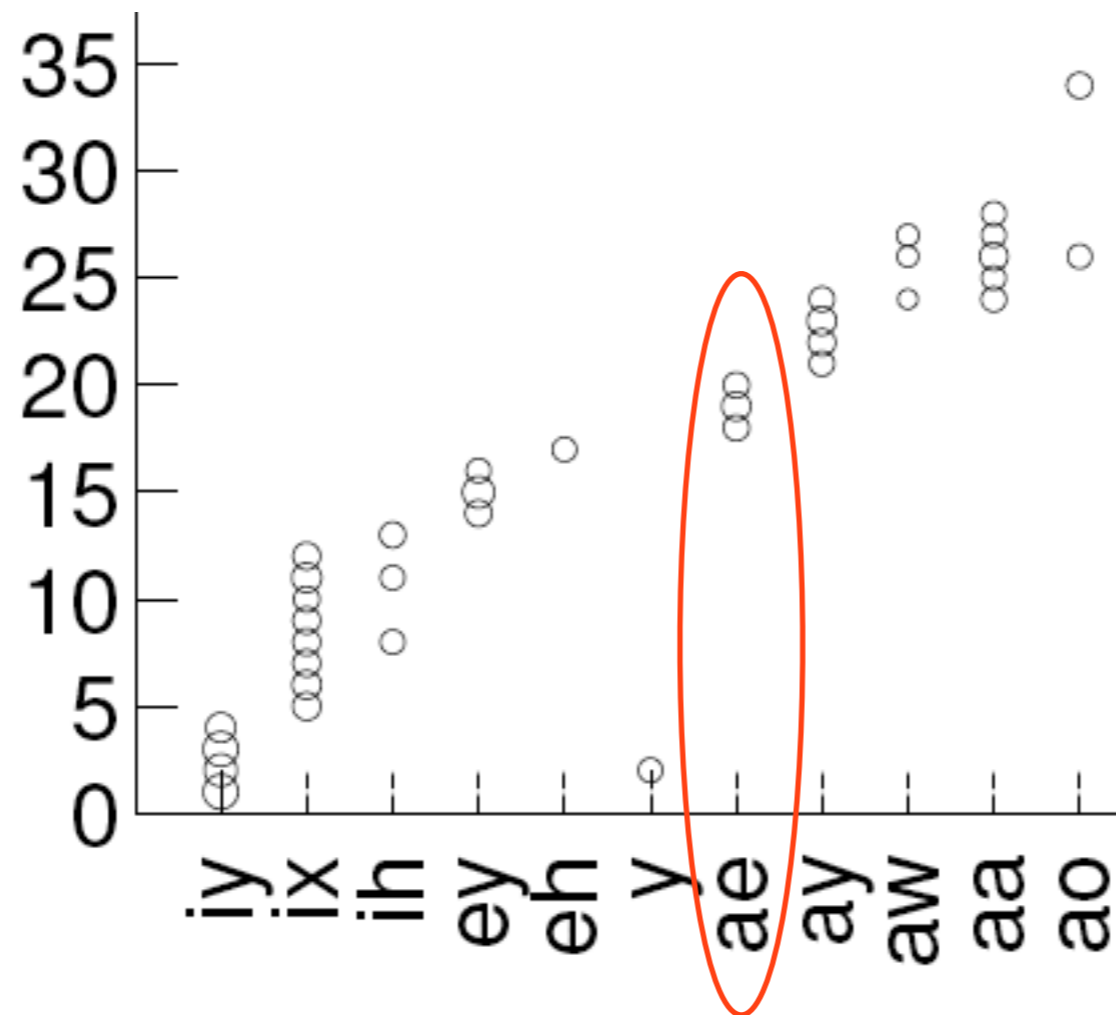- Three components in the model

  Adaptor grammars

  Noisy-channel model    Regularize the phonetic variations

  Phone discovery model

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

/k/ /ae/ /t/

49   58   32

Noisy-channel

/k/ /ae/ /t/          /k/ /ae/ /t/

49   58   32          49   26   32

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  – Substitution, deletion, insertion, and exact-match

/k/ /ae/ /t/

49   58   32

Noisy-channel

/k/ /ae/ /t/            /k/ /ae/ /t/

49   58   32            49   26   32

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

    – Substitution, deletion, insertion, and exact-match

/k/ /ae/ /t/

49  58  32

Noisy-channel

/k/ /ae/ /t/

49  58  32

/k/ /ae/ /t/

49  26  32

Three exact-match operations
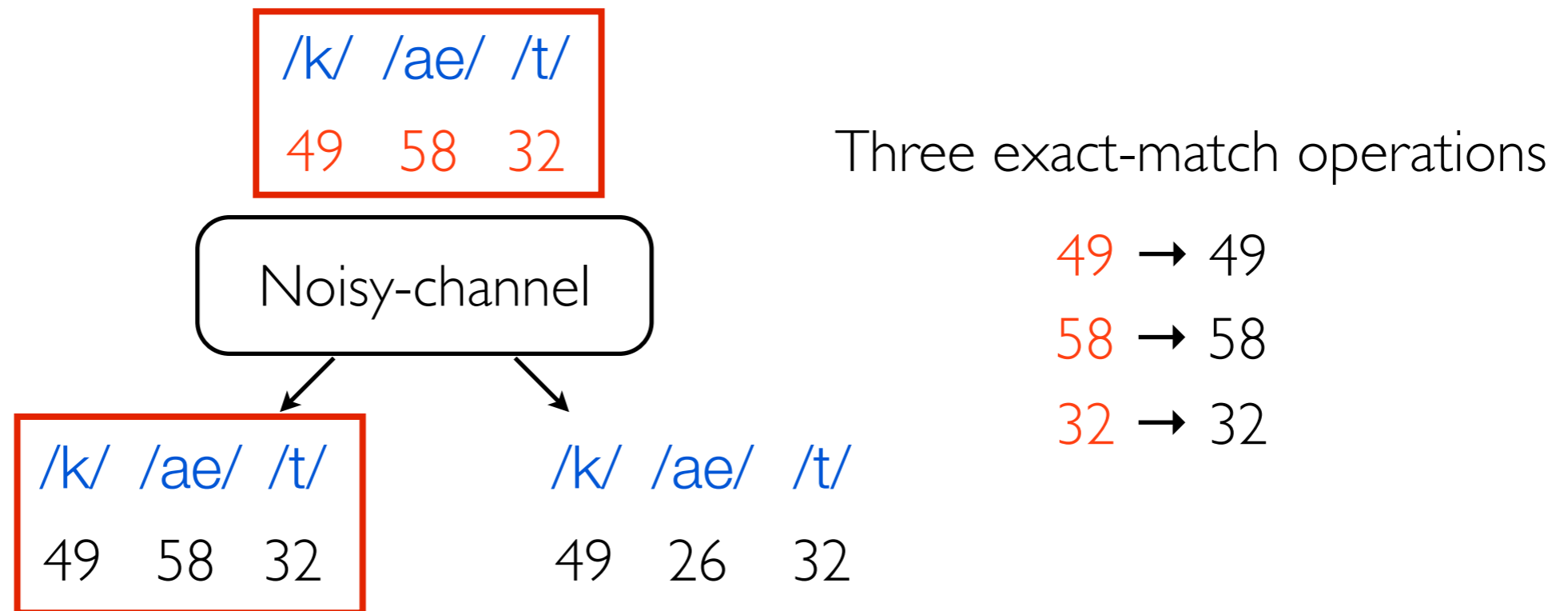
49 → 49

58 → 58

32 → 32

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  - Substitution, deletion, insertion, and exact-match

/k/  /ae/  /t/

49    58    32

Noisy-channel

/k/  /ae/  /t/

49    58    32

/k/  /ae/  /t/

49    26    32

exact-match  49 → 49

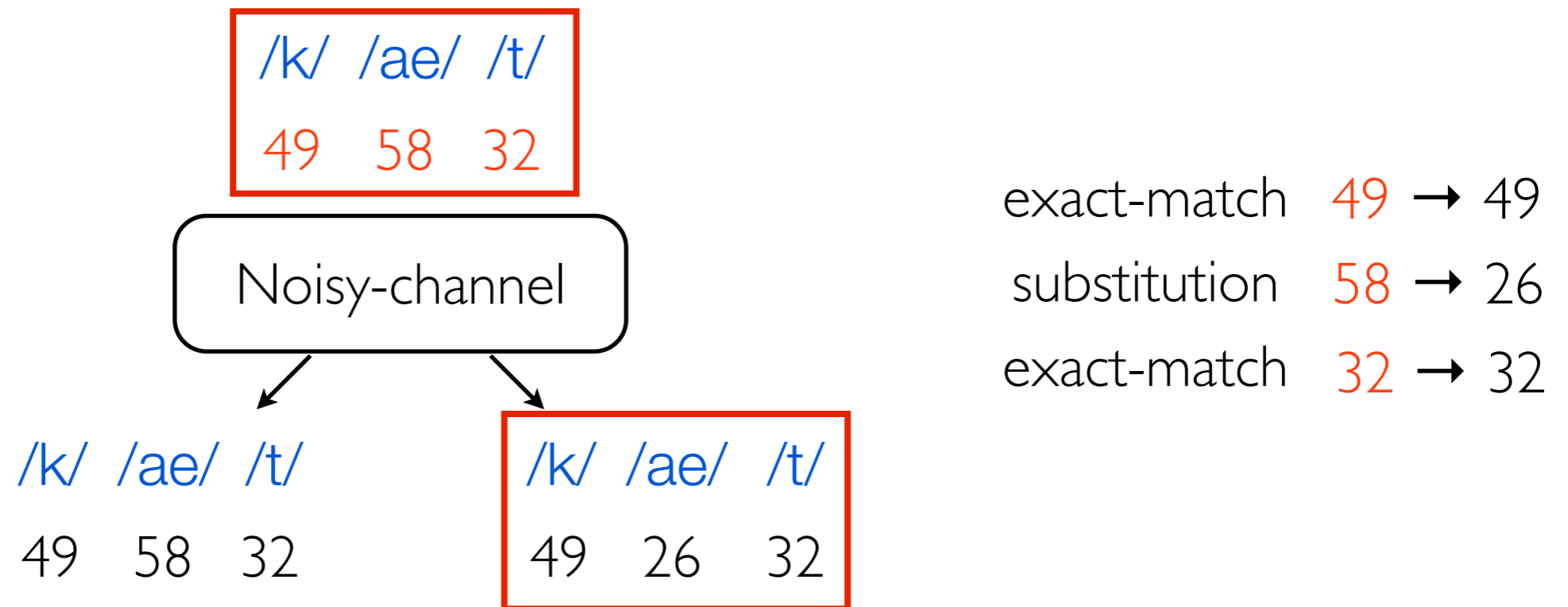substitution  58 → 26

exact-match  32 → 32

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

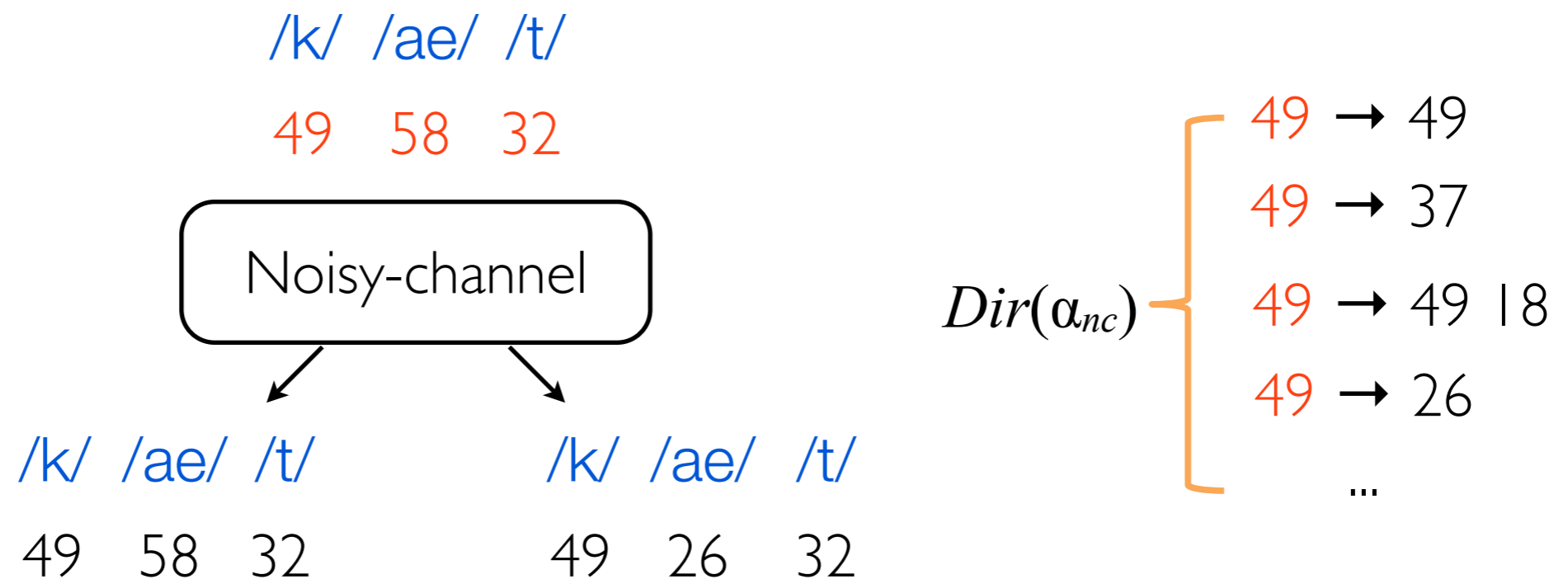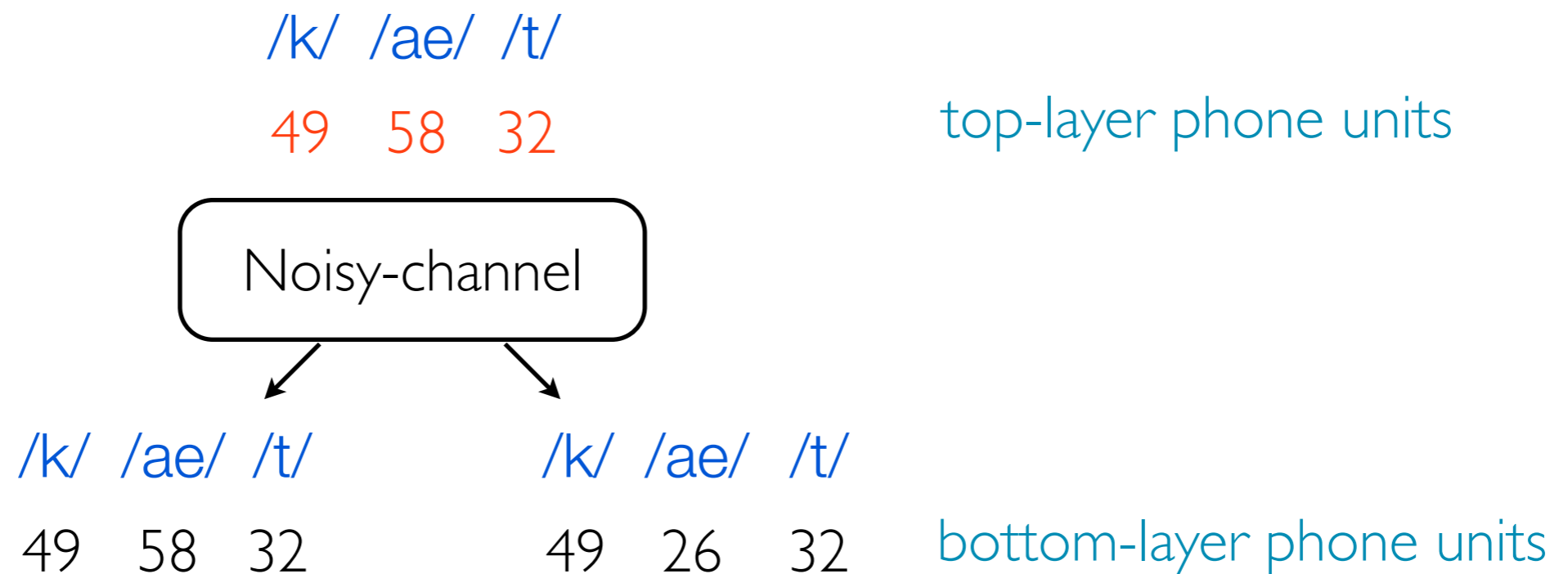  – Substitution, deletion, insertion, and exact-match

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  – Substitution, deletion, insertion, and exact-match

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phone discovery model

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  – To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phone discovery model

# Generative Process

- Generate a parse from adaptor grammars

| Adaptor grammars |
| Noisy-channel model |
| Phone discovery model |

# Generative Process

- Generate a parse from adaptor grammars



$d$    word and syllable structures

$u$    top-layer phone units

Adaptor grammars

Noisy-channel model

Phone discovery model

# Generative Process

- Generate phonetic variations

# Generative Process

- Generate phonetic variations

# Generative Process

- Generate phonetic variations

# Generative Process

- Generate phonetic variations

# Generative Process

- Generate phonetic variations

# Generative Process

- Generate speech data



d

Sen

Word

Syl        Syl

Phn  Phn

u    3    16    58    49

v    3    5    16    37    49
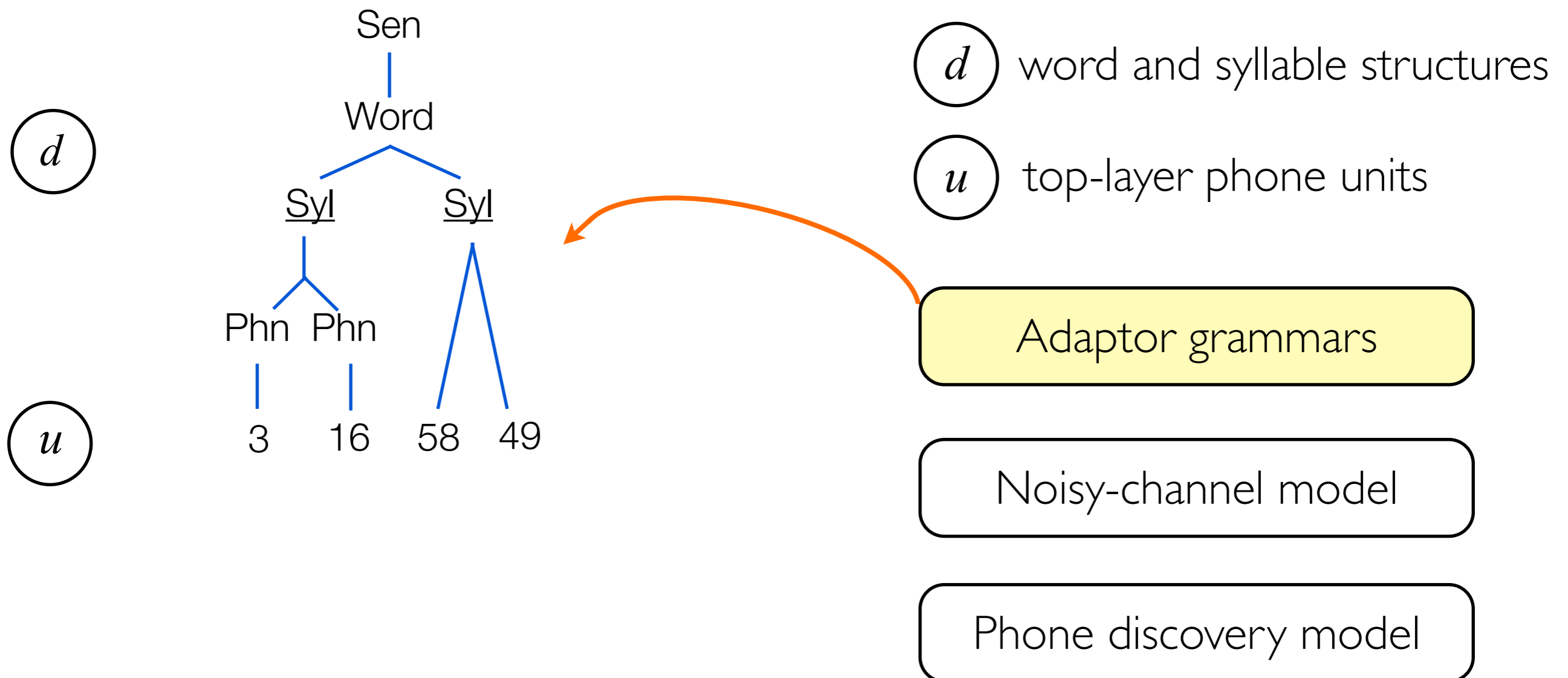
Adaptor grammars

Noisy-channel model

Phone discovery model

# Generative Process

- Generate speech data

# Generative Process

- Generate speech data

# Generative Process

- **Generate speech data**

# Generative Process

# Inference

# Inference

- Only speech data are observed

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phone discovery model

# Initialization

- Initialize $v$ and $b$ using the phonetic discovery model

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phone discovery model

# Initialization

- Initialize $v$ and $b$ using the phonetic discovery model

# Inference

- Given $v$ and $b$ sample $d$ and $u$

$d$

$u$

$v$

$x$

$b$

3　5　16　37　49

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given $v$ and $b$ sample $d$ and $u$

$d$

Metropolis-Hastings algorithm

$u$

$v$

$x$     3   5   16    37    49



$b$

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given $v$ and $b$ sample $d$ and $u$



Sen
Word
Syl    Syl
Phn Phn
3    16   58   49
3  5  16  37  49

$d$

$u$

$v$

$x$

$b$

Metropolis-Hastings algorithm

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given $d$ and $u$ resample $v$ and $b$

Sen
Word
Syl  Syl
Phn  Phn
3  16  58  49

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given $d$ and $u$ resample $v$ and $b$

$d$

Sen
|
Word
Syl      Syl
Phn  Phn
3    16    58    49

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given $d$ and $u$ resample $v$ and $b$

Sen

Word

Syl    Syl

Phn  Phn

3    16    58    49

3    16    37    49

Leverage higher level knowledge
to learn phonetic structures

Adaptor grammars

Noisy-channel model

Phone discovery model

$d$

$u$

$v$

$x$

$b$

# Inference

- Given $v$ and $b$ sample $d$ and $u$

$d$

$u$

$v$

$x$

$b$

3    16    37    49

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given ($v$) and ($b$) sample ($d$) and ($u$)

($d$)

($u$)
($v$)
($x$)
($b$)

3    16    37    49

Adaptor grammars

Noisy-channel model

Phone discovery model

# Inference

- Given (v) and (b) sample (d) and (u)

Alternate between the two steps

Sen

Word

Syl          Syl

Phn  Phn

3    17    58    49

3    16    37    49

(d)
(u)
(v)
(x)
(b)

Adaptor grammars

Noisy-channel model

Phone discovery model

# Experimental Setup

- **MIT Lecture Corpus**

  - The six lectures evaluated in [*Park and Glass, IEEE Trans. 2008*]

  - Each lecture contains ~1 hour of speech data by a single speaker

  - Each lecture contains a set of subject-specific keywords

- **Qualitative assessment**

  - Sentence and word parses

  - Analysis on the discovered hierarchical linguistic structures

- **Quantitative assessment**

  - Coverage of subject-specific keywords

  - (Word and phone segmentation)

# Parse of a Full Sentence

37  12  67  88  158  1  2  19  20  41  47  13  103  48  91  4  67  25  8  99  29  44  22  103  4  37  12  67

# Parse of a Full Sentence

# Parse of a Full Sentence



MIT's only occurs 3 times in the lecture

open and university almost always appear together in the lecture

# Word Parses

- Two instances of "collaboration"

# Word Parses

- Two instances of "collaboration"

| $u$ | 50 | 137 | 28 | 16 | 18 | 31 | 43 | 6 | 7 | 30 | 50 | 137 | 28 | 16 | 18 | 31 | 43 | 6 | 7 | 30 |
|-----|----|-----|----|----|----|----|----|---|---|----|----|-----|----|----|----|----|----|---|---|----|
| $v$ | 91 | 106 | 28 | 16 | 18 | 29 | 43 | 6 | 7 | 30 | 50 | 106 | 28 | 16 | 18 | 29 | 43 | 41 | 7 | 30 |

# Word Parses

- Two instances of "collaboration"

  – Noisy-channel model regularizes the bottom-layer phone units

collaboration

[50 137]   [28 16]   [18 31 43]   [6 7 30]

kcl  k      el ae      bcl ax r      ey sh en

$u$  | 50  137  28  16  18  31  43  6  7  30 |  | 50  137  28  16  18  31  43  6  7  30 |

$v$  91  106  28  16  18  29  43  6  7  30      50  106  28  16  18  29  43  41  7  30

# Word Parses

- Two instances of "collaboration"

  – Noisy-channel model regularizes the bottom-layer phone units

  – Highly reusable sub-word structures

collaboration

[50 137]   [28 16]   [18 31 43]   [6 7 30]

kcl  k       el ae       bcl ax r       ey sh en

$u$  | 50  137  28  16  18  31  43  6  7  30 | 50  137  28  16  18  31  43  6  7  30

$v$  | 91  106  28  16  18  29  43  6  7  30 | 50  106  28  16  18  29  43  41  7  30

# Structure Reuse

- Examples of reusing [6 7 30]



collaboration

[50 137]   [28 16]   [18 31 43]   **[6 7 30]**
 kcl  k      el ae      bcl ax r      ey sh en

reservation

[1 158]  [70 23]  [34 99]  **[6 7 30]**
 r  eh      z       er  v    ey sh en

innovation

[67]  [1 27]  [99]  **[6 7 30]**
 ih    n ax     v    ey sh en

globalization

[106 48]   [18 31]   [147 13]   **[6 7 30]**
gcl  g l ow  bcl ax l     ax z      ey sh en

foundation

[22 46 8]         **[6 7 30]**
 f  aw n dcl d  ey sh en

# Subject-specific Keywords

- **Term Frequency Inverse Document Frequency (TFIDF) scores**

  – The top 20 words for each lecture [*Park and Glass, IEEE Trans. 2008*]

- **Keyword examples**

  – From the seminar about the book "The world is flat" by Thomas Friedman

| | | | |
|---|---|---|---|
| 1. flat | 6. flattener | 11. airline | 16. huge |
| 2. globalization | 7. dollar | 12. thousand | 17. create |
| 3. collaboration | 8. China | 13. outsourcing | 18. convergence |
| 4. India | 9. southwest | 14. really | 19. connect |
| 5. era | 10. argue | 15. platform | 20. chapter |

# Coverage of Keywords

# Coverage of Keywords

# Coverage of Keywords



Legend:
- Park & Glass, 2008 (blue)
- Full model (red)
- No resampling of bottom-layer units (orange)

Synergies between the learning of words and phones

Y-axis: # discovered keywords (0, 5, 10, 15, 20)

X-axis: Lecture topics (Economics, Signal processing, Clustering, Speaker adaptation, Physics, Linear algebra)

# Coverage of Keywords

# Coverage of Keywords

# Conclusion

- Two models for discovering linguistic structures from speech

# Conclusion

- Two models for discovering linguistic structures from speech

Discovering phonetic inventory

/b/  /ax/  /n/  /ae/  /n/  /ax/

- DP mixture models with HMMs

  - Discovered phonetic units are highly correlated with standard phones

# Conclusion

- Two models for discovering linguistic structures from speech

**Discovering phonetic inventory**

/b/ /ax/ /n/ /ae/ /n/ /ax/



- DP mixture models with HMMs

  – Discovered phonetic units are highly correlated with standard phones

**Discovering hierarchical linguistic structures**

Word     banana

Syllable

Phone   /b/ /ax/ /n/ /ae/ /n/ /ax/

# Conclusion

- Two models for discovering linguistic structures from speech

**Discovering phonetic inventory**

/b/  /ax/  /n/  /ae/  /n/  /ax/



- DP mixture models with HMMs

  – Discovered phonetic units are highly correlated with standard phones

**Discovering hierarchical linguistic structures**

Word       banana

Syllable

Phone   /b/  /ax/  /n/  /ae/  /n/  /ax/



- Integrate adaptor grammars with the phone discovery model

# Conclusion

- Two models for discovering linguistic structures from speech

**Discovering phonetic inventory**

/b/  /ax/  /n/  /ae/  /n/  /ax/

- DP mixture models with HMMs

  – Discovered phonetic units are highly correlated with standard phones

**Discovering hierarchical linguistic structures**

Word              banana

Syllable

Phone  /b/  /ax/  /n/  /ae/  /n/  /ax/

- Integrate adaptor grammars with the phone discovery model

  – Noisy-channel model is critical for learning lexical units

# Conclusion

- Two models for discovering linguistic structures from speech

**Discovering phonetic inventory**

/b/ /ax/ /n/ /ae/ /n/ /ax/
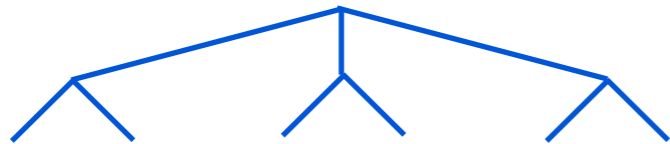


- DP mixture models with HMMs

  – Discovered phonetic units are highly correlated with standard phones

**Discovering hierarchical linguistic structures**

Word        banana

Syllable

Phone  /b/ /ax/ /n/ /ae/ /n/ /ax/



- Integrate adaptor grammars with the phone discovery model

  – Noisy-channel model is critical for learning lexical units

  – Synergies between word and phone learning

# Models and Applications

## Discovering phonetic inventory

*[Lee and Glass, ACL 2012]*

/b/  /ax/  /n/  /ae/  /n/  /ax/

## Discovering hierarchical linguistic structures

*[Lee, O'Donnell, and Glass, TACL 2015]*

Word      banana

Syllable

Phone   /b/  /ax/  /n/  /ae/  /n/  /ax/

# Models and Applications

## Discovering phonetic inventory

*[Lee and Glass, ACL 2012]*

/b/  /ax/  /n/  /ae/  /n/  /ax/

## One-shot learning spoken words

*[Lake\*, Lee\*, Glass, Tenenbaum, CogSci 2014]*

*\* share first authorship*

[k][ae][t]     [k][ae][t]     [d][ao][g]

## Discovering hierarchical linguistic structures

*[Lee, O'Donnell, and Glass, TACL 2015]*

Word     banana

Syllable

Phone   /b/  /ax/  /n/  /ae/  /n/  /ax/

# Models and Applications

## Discovering phonetic inventory

*[Lee and Glass, ACL 2012]*

/b/ /ax/ /n/ /ae/ /n/ /ax/



## One-shot learning spoken words

*[Lake\*, Lee\*, Glass, Tenenbaum, CogSci 2014]*

*\* share first authorship*

[k][ae][t]    [k][ae][t]    [d][ao][g]



## Discovering hierarchical linguistic structures

*[Lee, O'Donnell, and Glass, TACL 2015]*

Word            banana

Syllable

Phone    /b/ /ax/ /n/ /ae/ /n/ /ax/



## Pronunciation lexicon learning

*[Lee, Zhang, and Glass, EMNLP 2013]*

b a n a n a

/b/ /ax/ /n/ /ae/ /n/ /ax/

# Future Work

- Learning from more sensory data

  - Speech and visual streams

The doggie is sleeping

# Future Work

- Building spoken language systems based on discovered vocabulary

  - For low-resource languages or languages without a writing system

Thank you.
([kite.com](kite.com))

# Discovered Phone Units -- 300 utterances

- **43 phone units discovered from 300 TIMIT utterances**

    - Phone units are correlated with English broad phone classes

# Dirichlet Process (DP)

- Let's start with Dirichlet distribution

  - Dirichlet distribution is a distribution over the K-dim probability simplex

$$(\alpha_1, \alpha_2, \alpha_3) \quad \sum_{i=1}^{3} \alpha_i = 1$$

# Dirichlet Process (DP)

- Let's start with Dirichlet distribution

  - Dirichlet distribution is a distribution over the K-dim probability simplex

  - Assume we have 3 HMMs in the mixture

# Inference for HMM Parameters (θ)

- **HMM is used to model each phone**

  – Three states with only left-to-right and self transitions

  – Always start from the first state

  – A diagonal GMM is used for the emission distributions



$$\sum_{i=1}^{8} w_{1,i} N\left(u_{1,i}, \sigma_{1,i}^2\right) \quad \sum_{i=1}^{8} w_{2,i} N\left(u_{2,i}, \sigma_{2,i}^2\right) \quad \sum_{i=1}^{8} w_{3,i} N\left(u_{3,i}, \sigma_{3,i}^2\right)$$

# Inference for HMM Parameters (θ)

- HMM is used to model each phone

  - Three states with only left-to-right and self transitions

  - Always start from the first state

  - A diagonal GMM is used for the emission distributions

- Latent variables

  - Transition probabilities ($a$)

  - Mixture weights ($w$)

  - Means ($\mu$)

  - Variances ($\sigma^2$)



$$\sum_{i=1}^{8} w_{1,i} N\left(u_{1,i}, \sigma_{1,i}^2\right) \quad \sum_{i=1}^{8} w_{2,i} N\left(u_{2,i}, \sigma_{2,i}^2\right) \quad \sum_{i=1}^{8} w_{3,i} N\left(u_{3,i}, \sigma_{3,i}^2\right)$$

# Priors and Posteriors for HMM

- **Priors**

  - <u>Dirichlet distributions</u> for transition probabilities ($a$) and mixture weights ($w$)

  - <u>Normal-gamma distributions</u> for Gaussian parameters ($\mu$, $\sigma^2$)

# Priors and Posteriors for HMM

- Priors

  - <u>Dirichlet distributions</u> for transition probabilities ($a$) and mixture weights ($w$)

  - <u>Normal-gamma distributions</u> for Gaussian parameters ($\mu, \sigma^2$)

- Posteriors

  - Gather relevant counts from customer segments

# Priors and Posteriors for HMM

- **Priors**

  - <u>Dirichlet distributions</u> for transition probabilities ($a$) and mixture weights ($w$)

  - <u>Normal-gamma distributions</u> for Gaussian parameters ($\mu$, $\sigma^2$)

- **Posteriors**

  - Gather relevant counts from customer segments

  - Update prior distributions

  - Sample new values for the latent variables

# Dirichlet Process (DP)

- **Conceptually**

  – Dirichlet process can be viewed as an infinite case of Dirichlet distribution



- **Unknown # of HMMs**

  – Assume there are infinite number of HMMs first

  – Infer the finite number of HMM are needed to explain the finite data

  – By integrating $\beta$ during inference, DP provides a nice math format to find the #

# PCFG Review

- A PCFG is a quintuple $(N, T, S, R, \{\pi^q\}_{q \in N})$

- $N$ : a finite set of <u>nonterminal</u> symbols

- $T$ : a finite set of <u>terminal</u> symbols

  - $N \cap T = \varnothing$

- $S$ : start symbol

  - $S \in N$

- $R$ : production rules

  - $R = \{N \rightarrow (N \cup T)^*\}$

- $\pi^q$ : rule probabilities

  - $q \in N$

PCFG

| 0.5 | Sen | $\longrightarrow$ | Word Word |
| 0.5 | Sen | $\longrightarrow$ | Word |
| 0.7 | Word | $\longrightarrow$ | Syl Syl |
| 0.3 | Word | $\longrightarrow$ | Syl |
| 0.6 | Syl | $\longrightarrow$ | Phn Phn |
| 0.4 | Syl | $\longrightarrow$ | Phn |
| 0.1 | Phn | $\longrightarrow$ | /ax/ |
| 0.05 | Phn | $\longrightarrow$ | /n/ |
| | | | ... |

# Acoustic Landmarks

- Naively, every frame can be a phone boundary

  - In fact, some frames are more likely to be boundaries and some are less likely

  - Compute landmarks [Glass et al. 2003] and only do inference on landmarks

  - A language-independent method



- Disadvantage

  - Put an upper bound on recall rate

- Advantage

  - Reduce inference load

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

    - 3696 utterances for discovering phone units

    - Compute posterior-grams on the HMM states of the discovered phone units

$$x : \text{a single frame of feature vector}$$

$$State_{i,j} : \text{the j-th state of the i-th HMM}$$

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

$x$ : a single frame of feature vector

$State_{i,j}$ : the j-th state of the i-th HMM

$$\text{posterior-gram(x)} = \left[ \frac{p(State_{i,j} \mid x)}{\sum_{i=1}^{K} \sum_{j=1}^{3} p(State_{i,j} \mid x)} \right] \text{ for } 1 \leq i \leq K \text{ and } 1 \leq j \leq 3$$

$K$ : the total number of HMMs

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

  - Apply dynamic time warping to keyword detection [Zhang et al, 2009]

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

  - Apply dynamic time warping to keyword detection [Zhang et al, 2009]

  - 10 selected keywords

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

  - Apply dynamic time warping to keyword detection [Zhang et al, 2009]

  - 10 selected keywords

P@N: the average precision of top N hits

| P@N | EER |
|-----|-----|
|     |     |

# Spoken Term Detection

- Given a spoken query (w), find all spoken documents that contain w

    - 3696 utterances for discovering phone units

    - Compute posterior-grams on the HMM states of the discovered phone units

    - Apply dynamic time warping to keyword detection [Zhang et al, 2009]

    - 10 selected keywords

P@N: the average precision of top N hits

|  | P@N | EER |
|---|---|---|
| English Monophone (Supervised) | 74 | 11.8 |
| Thai Monophone Model (Supervised) | 56.6 | 14.9 |
| Our model | 63 | 16.9 |

# Spoken Term Detection

- **Given a spoken query (w), find all spoken documents that contain w**

  - 3696 utterances for discovering phone units

  - Compute posterior-grams on the HMM states of the discovered phone units

  - Apply dynamic time warping to keyword detection [Zhang et al, 2009]

  - 10 selected keywords

P@N: the average precision of top N hits

|  | P@N | EER |
|---|---|---|
| English Monophone (Supervised) | 74 | 11.8 |
| Thai Monophone Model (Supervised) | 56.6 | 14.9 |
| Our model | 63 | 16.9 |
| Zhang 2009 (GMM) (Unsupervised) | 52.5 | 16.4 |
| Zhang 2012 (DBM) (Unsupervised) | 51.1 | 14.7 |

# Unknown Number of HMMs
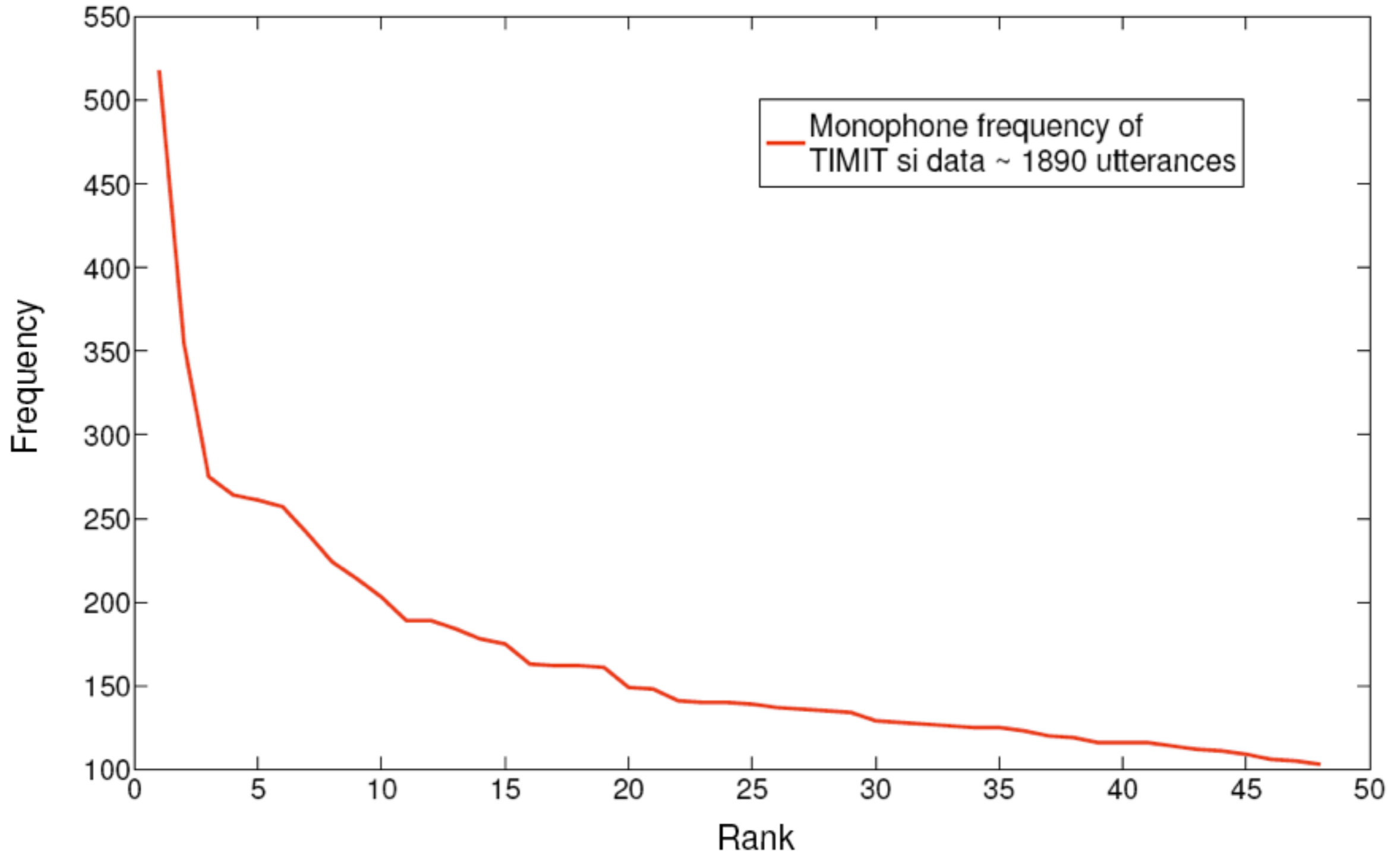
- An unknown set of phone units

# Unknown Number of HMMs

- An unknown set of phone units

  - Impose a Dirichlet Process prior to infer the number of phones

# Unknown Number of HMMs

- An unknown set of phone units

    - Impose a Dirichlet Process prior to infer the number of phones

- Is Dirichlet process (DP) a proper prior for this task?

    - Does phone frequency inherit power law?

# Phone Frequency -- Monophone
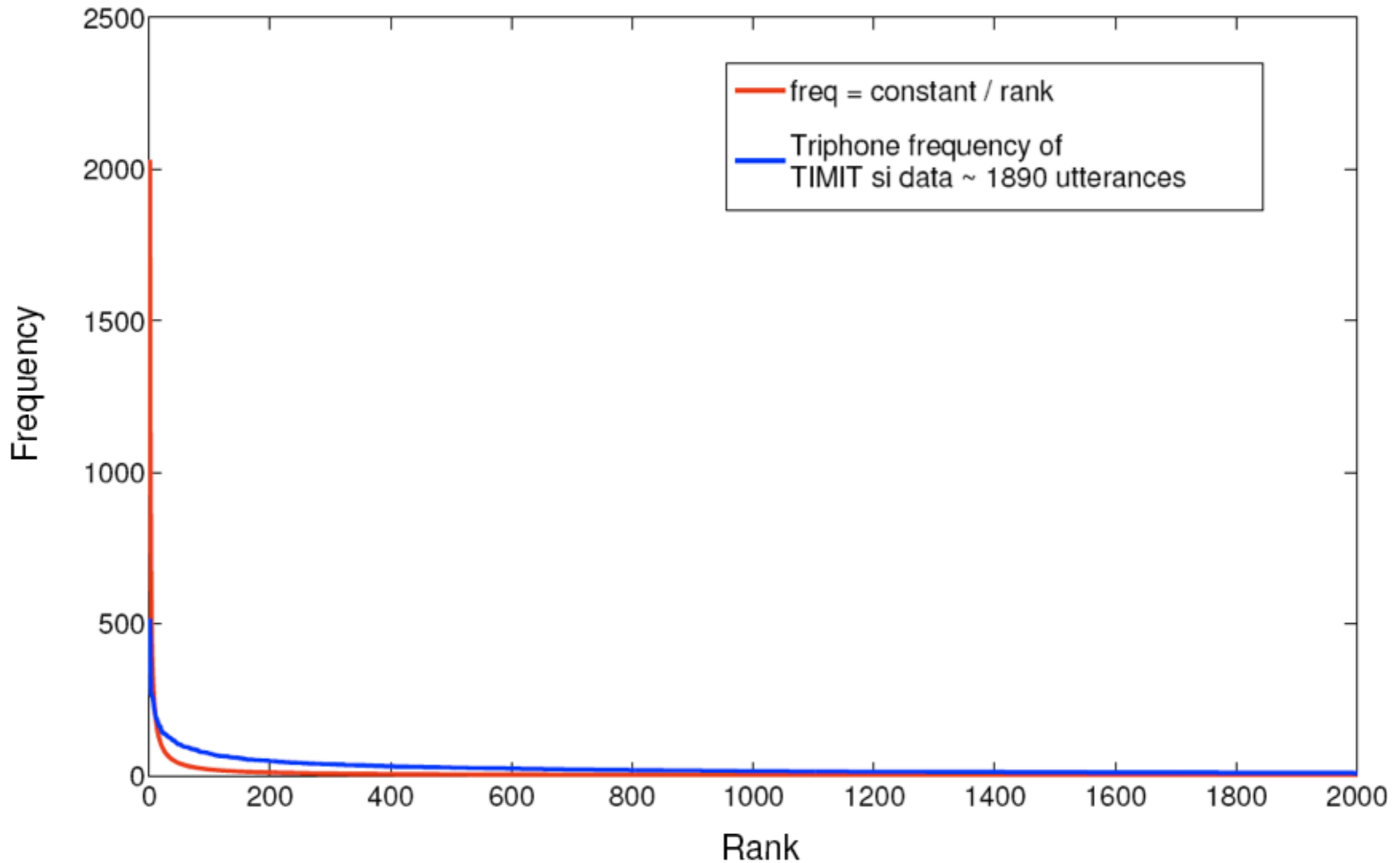
# Phone Frequency -- Triphone

# Unknown Number of HMMs

- An unknown set of phone units

  - Impose a Dirichlet Process prior to infer the number of phones

- Is Dirichlet process (DP) a proper prior for this task?

  - Does phone frequency inherit power law?

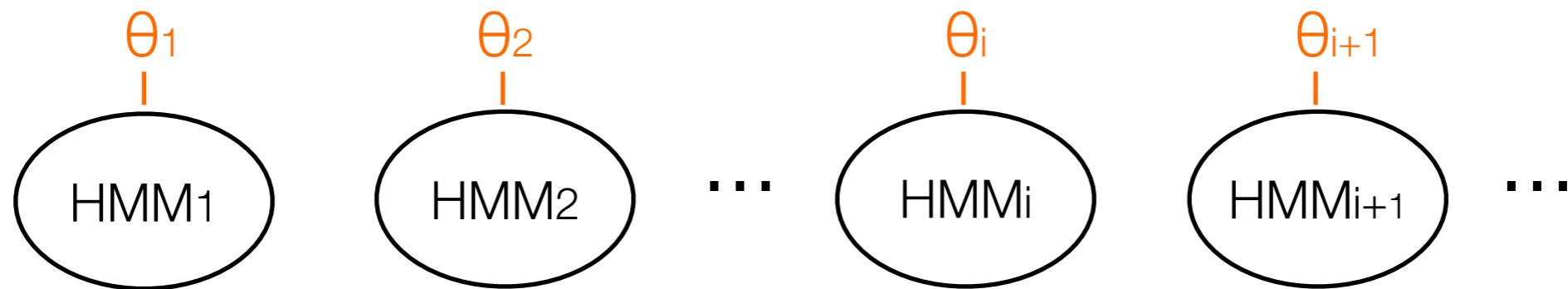# Unknown Number of HMMs

- An unknown set of phone units

    – Impose a Dirichlet Process prior to infer the number of phones

- Is Dirichlet process (DP) a proper prior for this task?

    – Does phone frequency inherit power law?

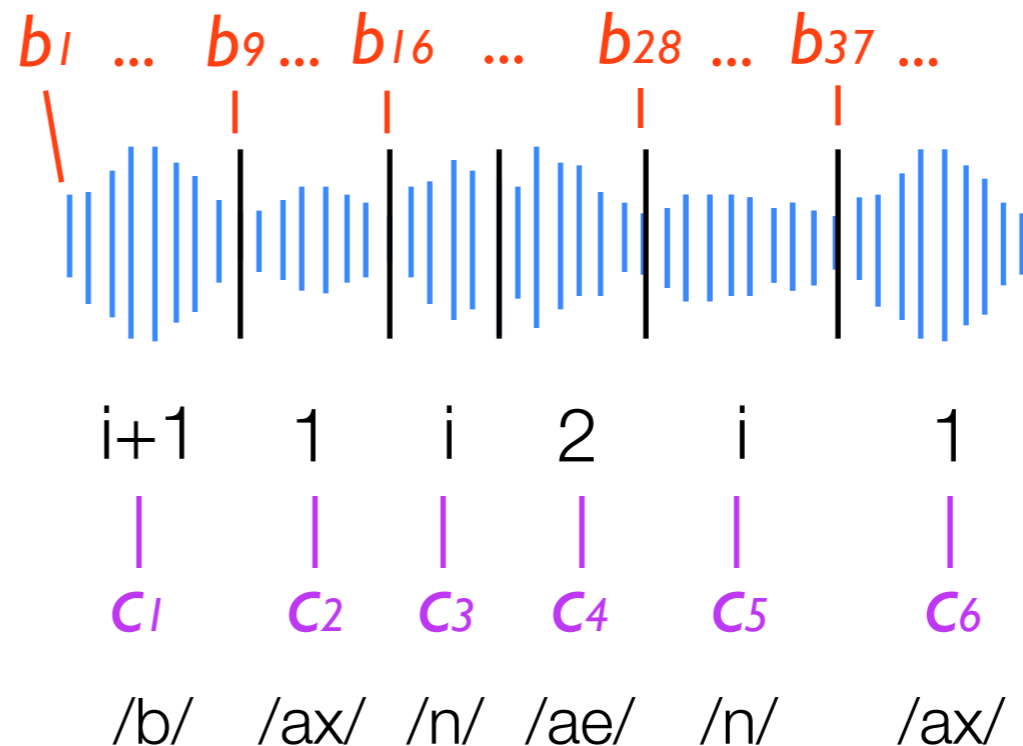    – DP should be a reasonable prior to start with

# Generative Story

- A simple explanation of how a spoken utterance is generated



- Main latent variables

  - Phone boundaries (*b*)

  - Phone labels (*c*)

  - HMM parameters (θ)

  - # of HMMs (phones)

    Dirichlet Process

# Language Acquisition Modeling

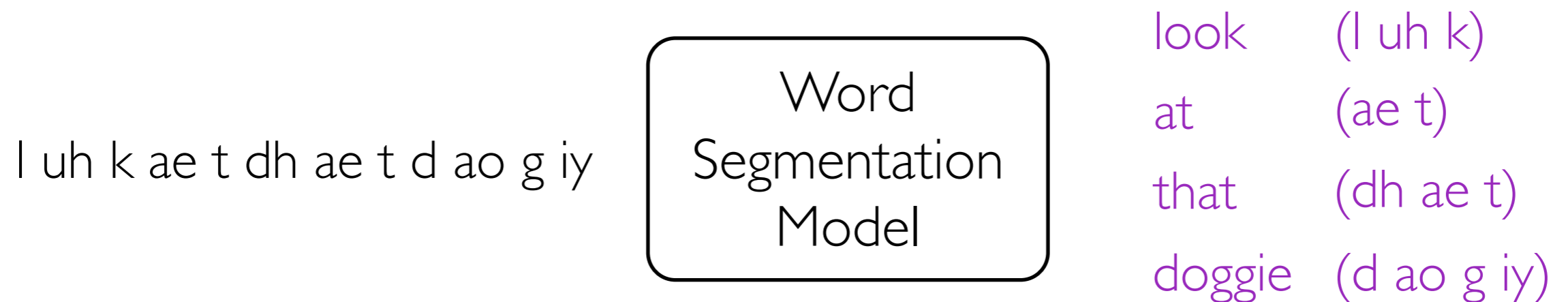- Previous work relies on highly pre-processed input data

# Language Acquisition Modeling

- Previous work relies on highly pre-processed input data

l uh k ae t dh ae t d ao g iy

| Word Segmentation Model |

look     (l uh k)
at       (ae t)
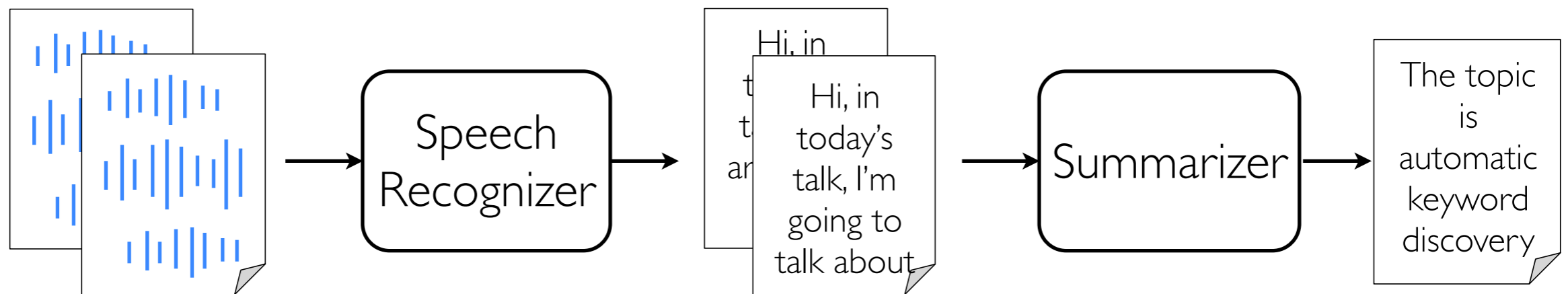that     (dh ae t)
doggie   (d ao g iy)

Other tasks such as phonetic unit learning are ignored

- Ground language acquisition modeling in real sensory data

- Ultimately allow machines to acquire a language like humans
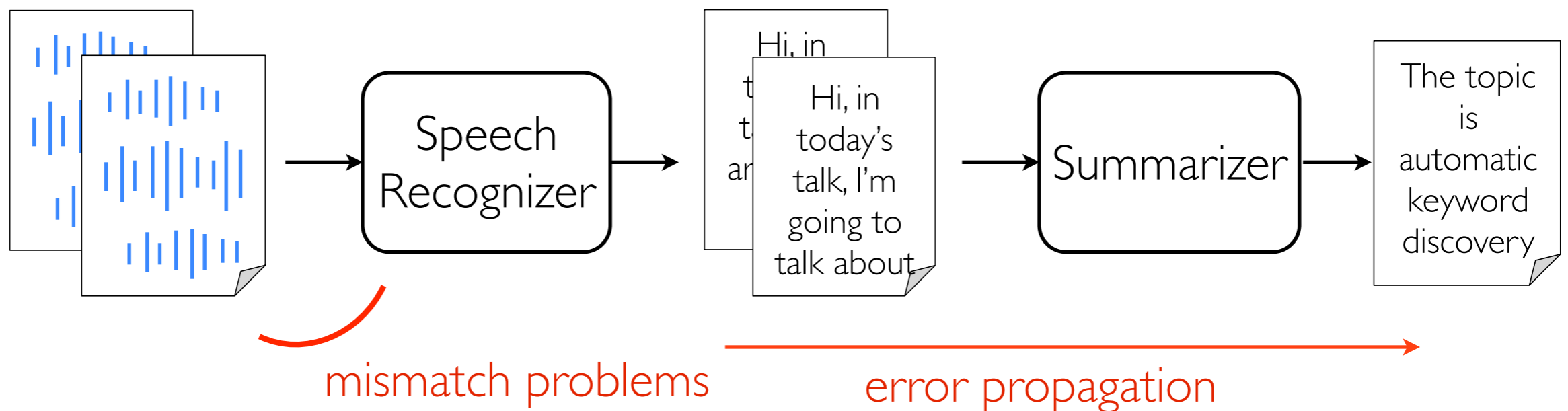
# Discovering Structures Beyond Phones

- Useful for representing out-of-vocabulary words
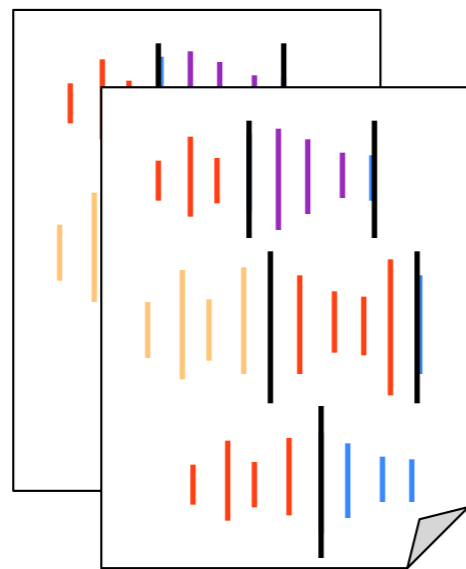
- Spoken document summarization

# Discovering Structures Beyond Phones

- Useful for representing out-of-vocabulary words

- Spoken document summarization
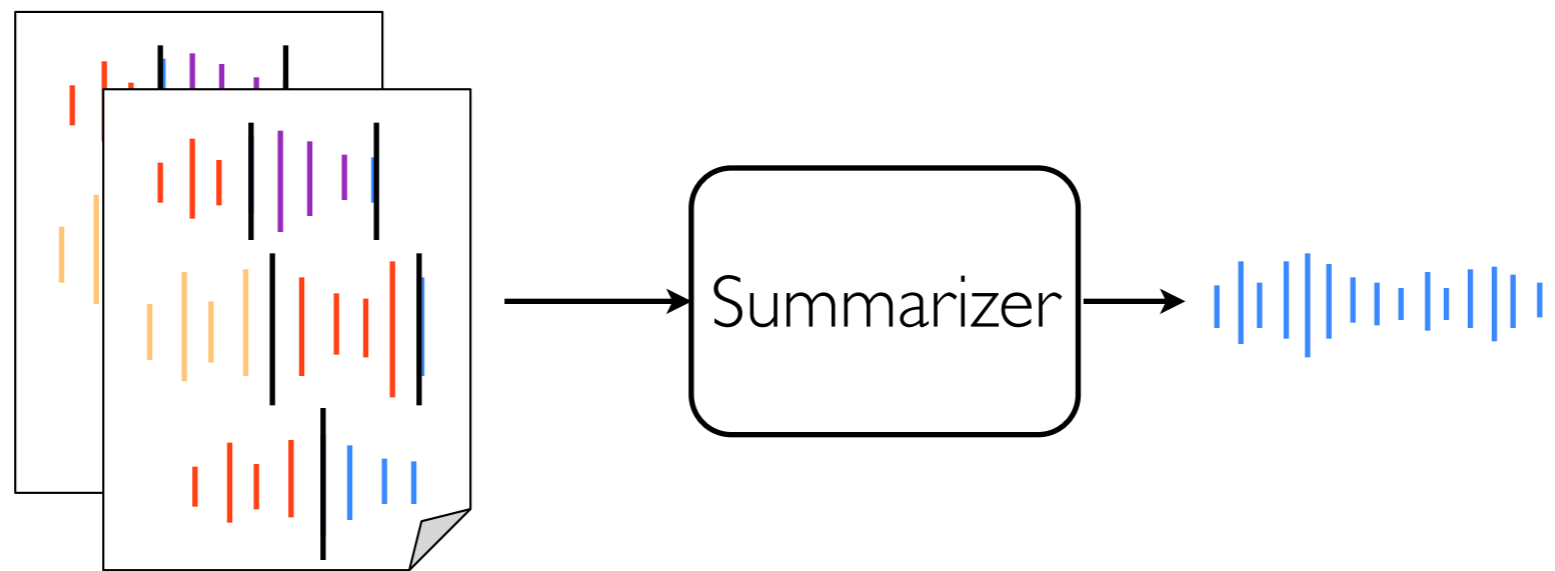
# Discovering Structures Beyond Phones

- Sub-word units are useful for representing out-of-vocabulary words

- Unsupervised word discovery
  - Automatic spoken document summarization without speech recognition



latent word structures

# Discovering Structures Beyond Phones

- **Sub-word units are useful for representing out-of-vocabulary words**

- **Unsupervised word discovery**

  – Automatic spoken document summarization without speech recognition



- **Connection to Cognitive Science (CogSci)**

  – Computational models for learning from speech are of great interests in CogSci

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

- Three components in the model

  Adaptor grammars

  Noisy-channel model

  Phonetic discovery model

  Discover hierarchical linguistic structures
  (Words, Syllables etc)

# Model Overview

- Integrate adaptor grammars and the phone discovery model

  - To discover rich linguistic structures from speech

- Three components in the model

Adaptor grammars

Noisy-channel model

Phonetic discovery model          Discover the phonetic units from acoustic data

# Model Overview

- Integrate adaptor grammars and the phone discovery model

    - To discover rich linguistic structures from speech

- Three components in the model
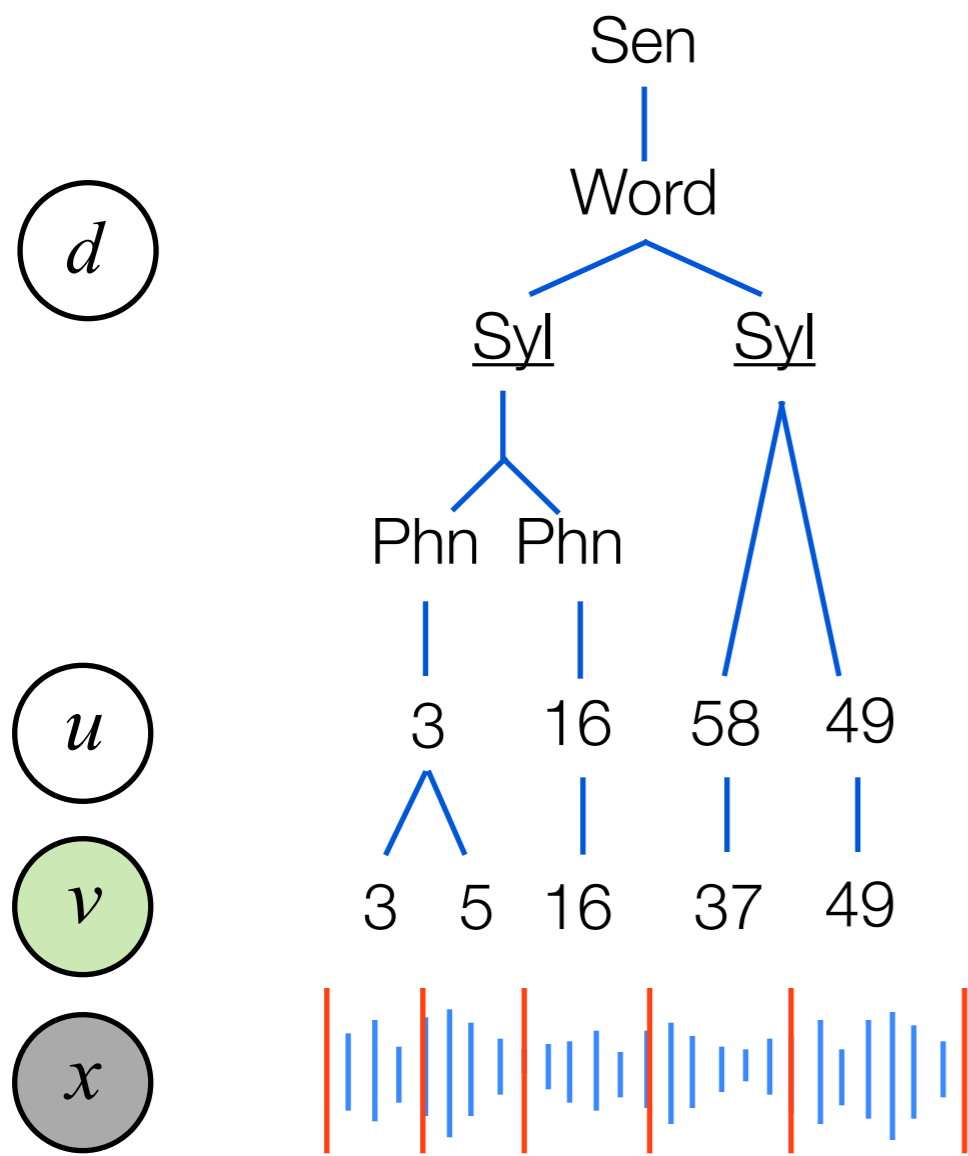
Adaptor grammars

Noisy-channel model    Bridges the other two components

Phonetic discovery model

# Inference

- Given $d$ and $u$



Adaptor grammars

Noisy-channel model

Phonetic discovery model

# Inference

- Given $d$ and $u$ resample $v$ and $b$



Sen

Word

Syl     Syl

Phn   Phn

3    16    58    49

3   5   16   37   49

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phonetic discovery model

# Initialization

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

Phonetic discovery model

# Initialization

$d$

$u$

$v$

$x$

$b$

Adaptor grammars

Noisy-channel model

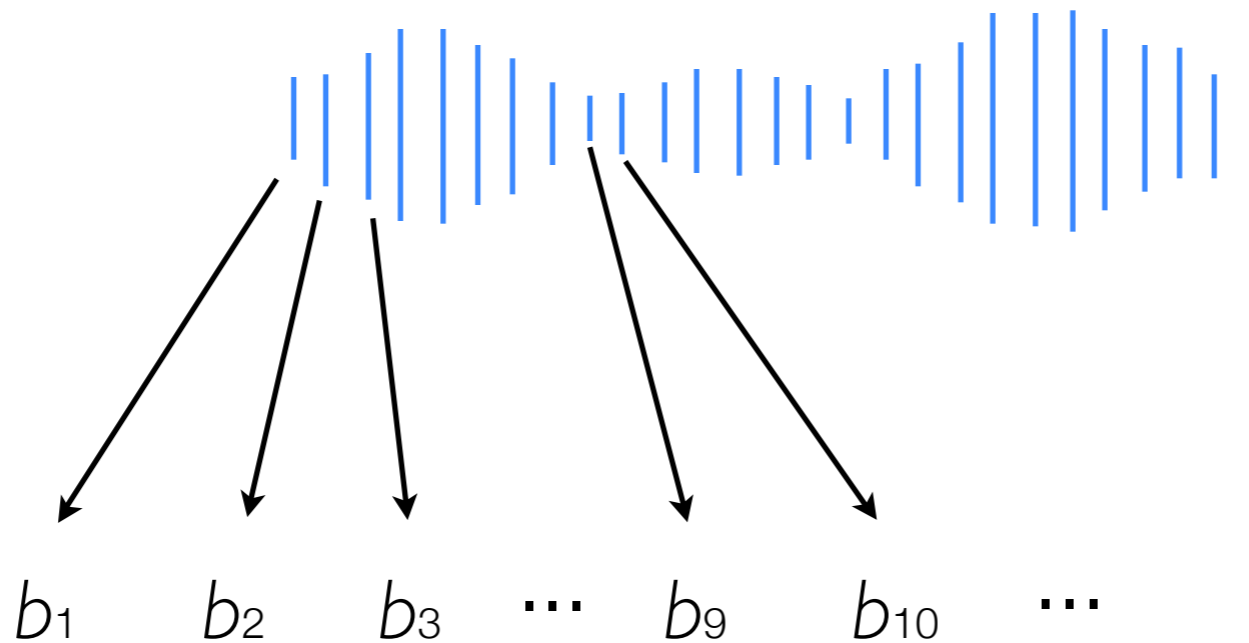Phonetic discovery model

# Initialization

- Initialize $v$ and $b$ using the phonetic discovery model

# Inference on Phone Boundaries ($b$)
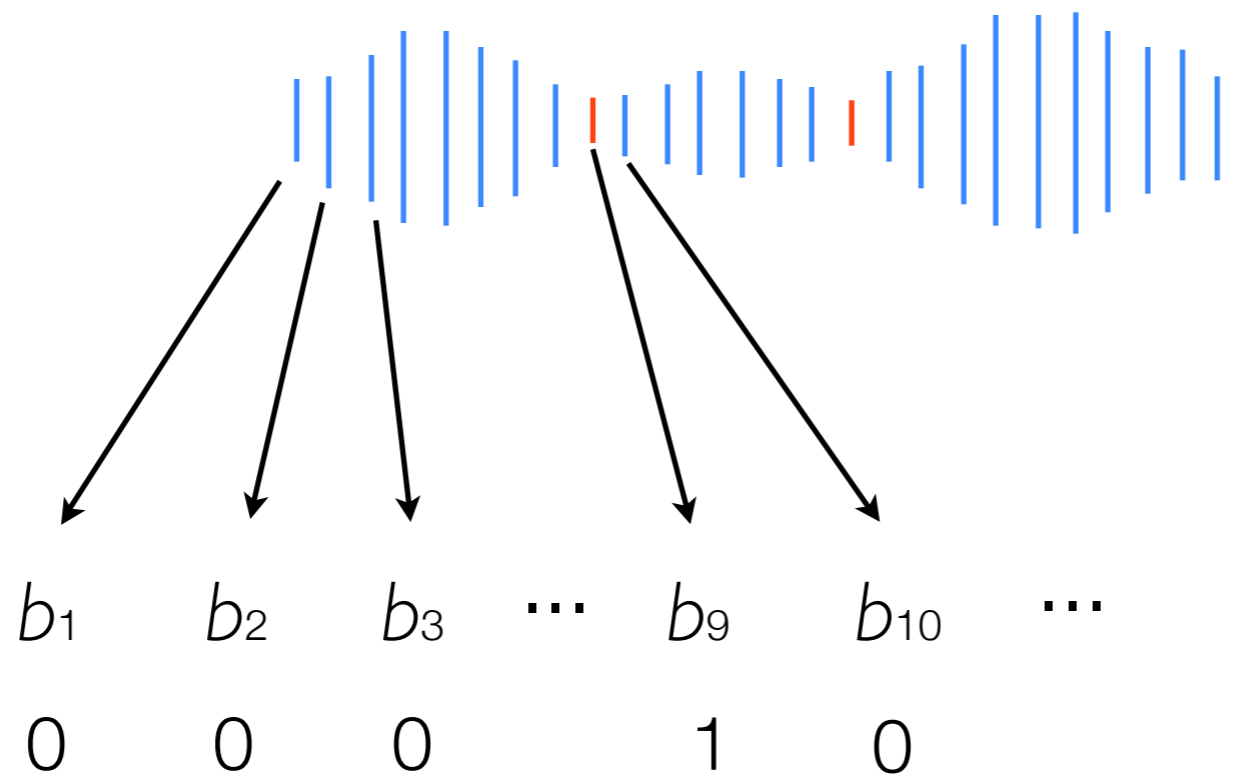
- **Boundary variables**

  - A priori, every frame can be a phone boundary

# Inference on Phone Boundaries ($b$)

- **Boundary variables**

  - A priori, every frame can be a phone boundary

  - Boundary variables take binary values

$$
\begin{array}{cccccc}
b_1 & b_2 & b_3 & \cdots & b_9 & b_{10} & \cdots \\
0 & 0 & 0 & & 1 & 0
\end{array}
$$

# Prior and Posterior for Phone Boundaries

- Prior

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$
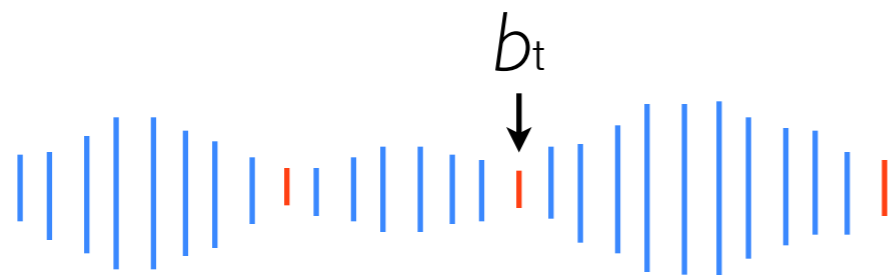
# Prior and Posterior for Phone Boundaries

- Prior

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

- Posterior: examine one boundary variable ($b_t$) at a time

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes

# Prior and Posterior for Phone Boundaries

- ## Prior

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

- ## Posterior: examine one boundary variable ($b_t$) at a time

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes
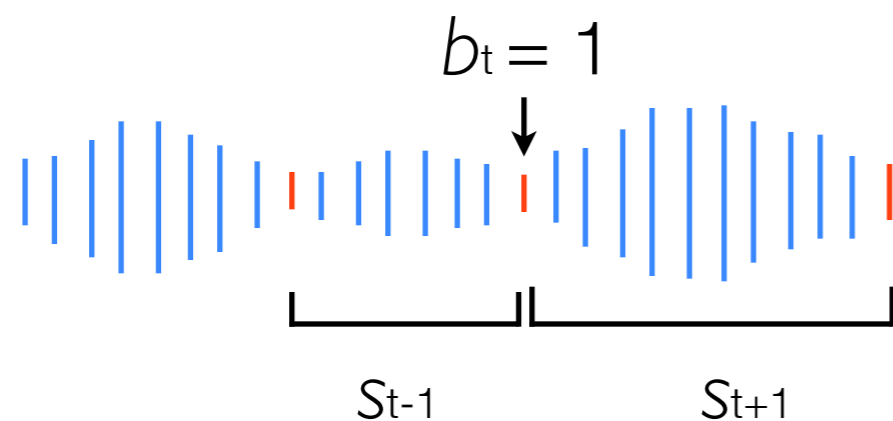
# Prior and Posterior for Phone Boundaries

- **Prior**

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

- **Posterior: examine one boundary variable ($b_t$) at a time**

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes



$$p(b_t = 1 \mid \cdots) \propto$$
$$p(b_t = 1) p(s_{t-1} \mid c^-, \underline{\theta}) p(s_{t+1} \mid c^-, \underline{\theta})$$

$c^-$ : cluster labels of all other segments

$\underline{\theta}$ : the set of HMMs
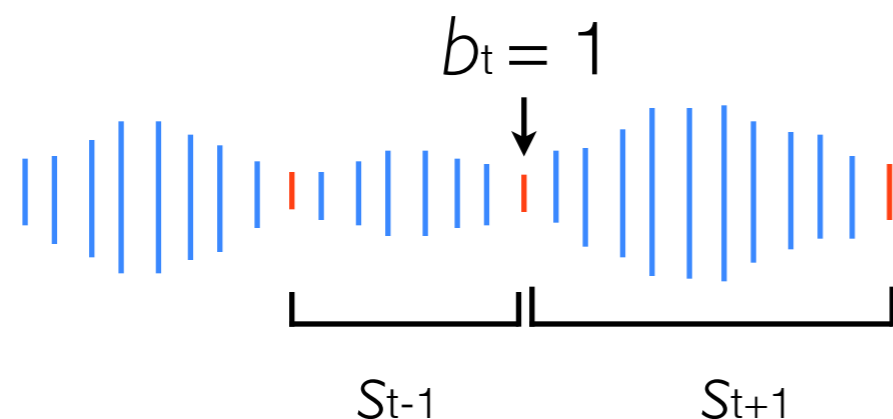
# Prior and Posterior for Phone Boundaries

- Prior

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

- **Posterior: examine one boundary variable ($b_t$) at a time**

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes



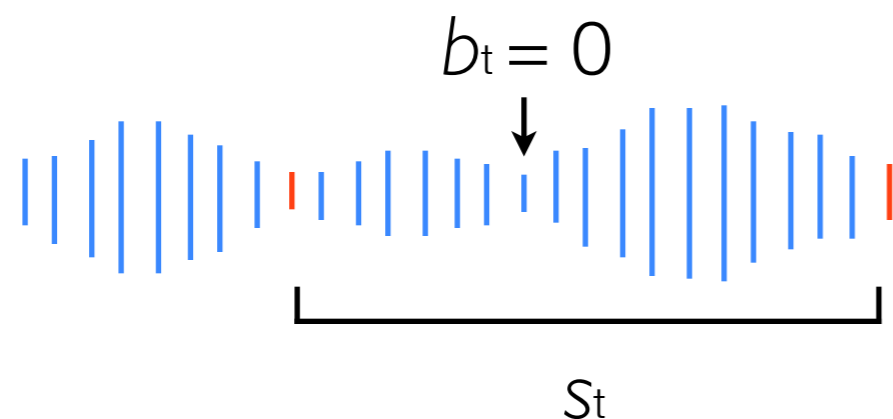$b_t = 0$

$s_t$

# Prior and Posterior for Phone Boundaries

- Prior

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

- **Posterior: examine one boundary variable ($b_t$) at a time**

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes

$b_t = 0$



$s_t$

$$p(b_t = 0 \mid \cdots) \propto$$
$$p(b_t = 0)p(s_t \mid c^-, \underline{\theta})$$
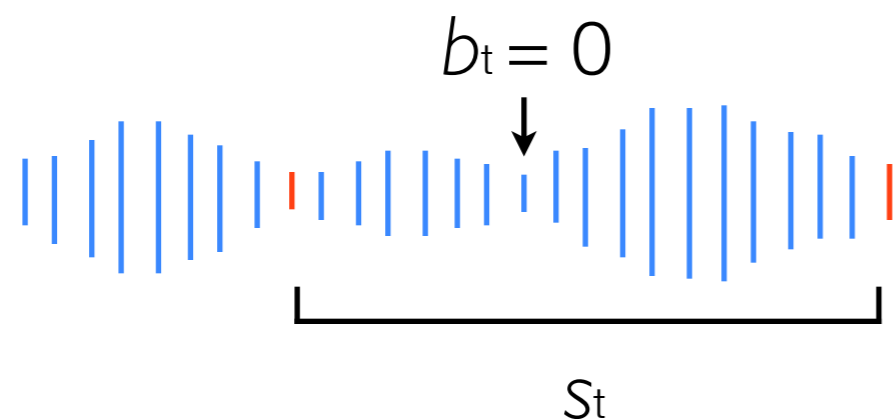
# Prior and Posterior for Phone Boundaries

- **Prior**

  - Fixed prior probabilities $p(b_t = 1) = \alpha_b$ and $p(b_t = 0) = 1 - \alpha_b$

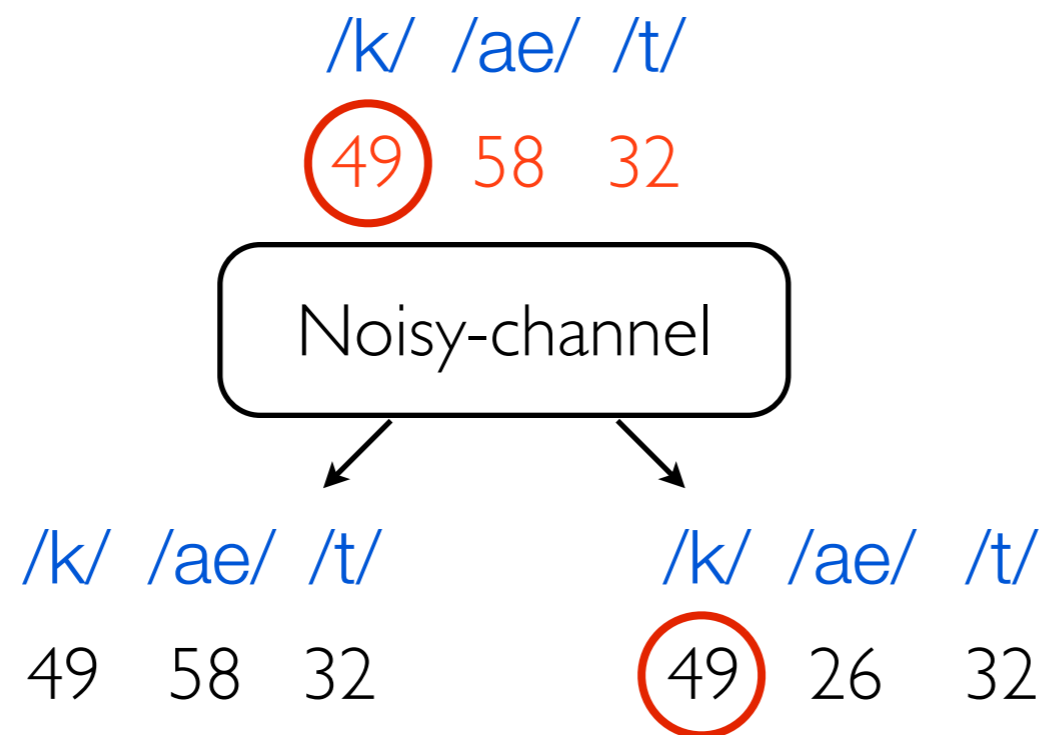- **Posterior: examine one boundary variable ($b_t$) at a time**

  - Fix the current values of other boundary variables

  - Consider both 0 and 1 for $b_t$ and the respective segmentation outcomes

Generate a sample for $b_t$

$$p(b_t = 1 \mid \cdots) \propto$$
$$p(b_t = 1)\, p(s_{t-1} \mid c^-, \underline{\theta})\, p(s_{t+1} \mid c^-, \underline{\theta})$$

$$p(b_t = 0 \mid \cdots) \propto$$
$$p(b_t = 0)\, p(s_t \mid c^-, \underline{\theta})$$

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  - Substitution, deletion, insertion, and exact-match



exact-match   49 → 49
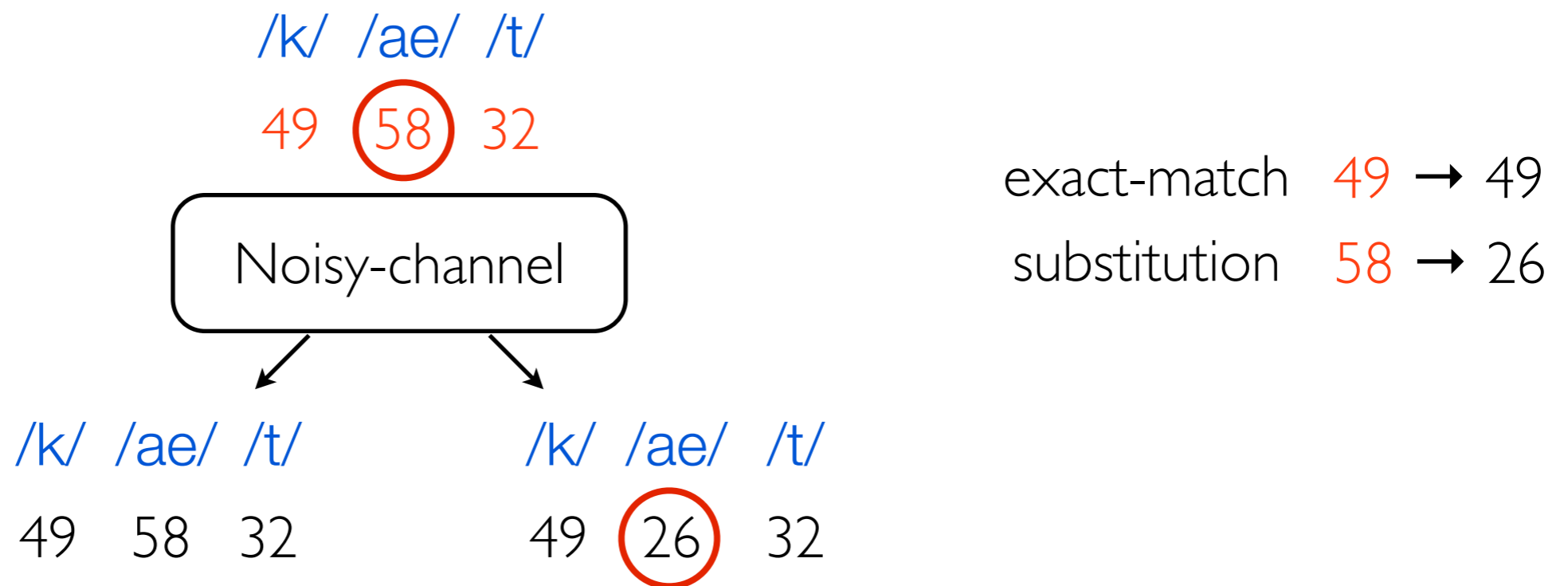
# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  – Substitution, deletion, insertion, and exact-match

/k/ /ae/ /t/

49 58 32

Noisy-channel

/k/ /ae/ /t/          /k/ /ae/ /t/

49 58 32          49 26 32

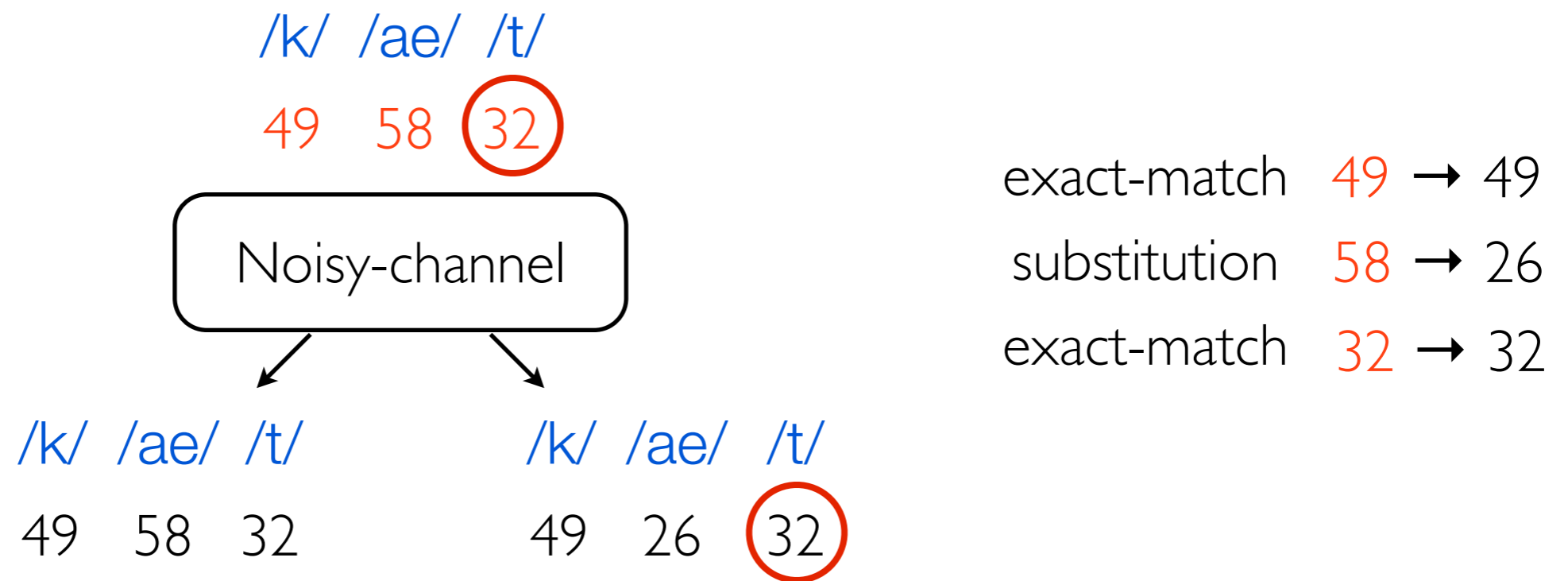exact-match 49 ➝ 49

substitution 58 ➝ 26

# Noisy-channel Model

- Assume the phonetic variations are outcomes of a noisy-channel

- Formulate the noisy-channel model as a set of edit operations

  - Substitution, deletion, insertion, and exact-match



exact-match  49 → 49

substitution  58 → 26

exact-match  32 → 32

# Acknowledgement

- **Thesis advisor**
  - Jim Glass

- **Thesis committee**
  - Regina Barzilay and Victor Zue

- **Collaborators**
  - Tim O'Donnell
  - Brenden Lake
  - Matt Johnson      – Yu Zhang
  - Lee Hetherington  – Ian McGraw
  - Stephanie Seneff   – Oded Ghitza

- **SLS members**
  - Marcia    – Tuka      – Xue      – Daniel
  - Scott     – Carrie    – David    – Michael
  - Najim     – Ekapol    – Mandy    – Stephen
  - Patrick   – Jennifer  – Ann      – Yu
  - Chengjie

- **Previous officemates**
  - Hung-an    – Paul
  - Yaodong    – Yuan

- **Friends and family**