# Bayesian Models of Lexicon Learning

Timothy J. O'Donnell
McGill Linguistics

# Language Learning and the Lexicon

Language learning consists mostly (entirely?) of learning the lexicon (words, morphemes, etc.)?

# The Problem of Lexical Uncertainty

How do learners identify the lexical units in their language?

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| | -ness |
| | -th |
| | -ity |

# The Problem of Lexical Uncertainty

Nouns (`N`)

Attach to adjectives (`Adj`)

Mean abstract quality or state

*goodness, cheapness, forgiveness, circuitousness, grandness, orderliness, business, goodness, …*

| | Suffix |
|---|---|
| | -ness |
| | -th |
| | -ity |

# The Problem of Lexical Uncertainty

Nouns (N)

Attach to adjectives (Adj)

Mean abstract quality or state

| | Suffix |
|---|---|
| *goodness, cheapness, forgiveness, circuitousness, grandness, orderliness, business, goodness, …* | -ness |
| *truth, warmth, width, depth, filth, sloth, strength, death, dearth, wealth, length, youth, …* | -th |
| | -ity |

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Nouns (N) | |
| Attach to adjectives (Adj) | |
| Mean abstract quality or state | |
| *goodness, cheapness, forgiveness, circuitousness, grandness, orderliness, business, goodness, …* | -ness |
| *truth, warmth, width, depth, filth, sloth, strength, death, dearth, wealth, length, youth, …* | -th |
| *verticality, tractability, severity, seniority, inanity, electricity, parity, scarcity, reality, …* | -ity |

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | **-ness** |

*pine-scented*          *pine-scentedness*

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | -ness |

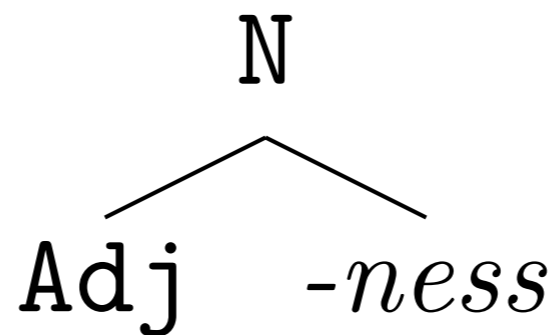*goodness, cheapness, forgiveness, circuitousness, grandness, orderliness, pretentiousness, business, goodness, greenness, …*

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | -ness |

*good**ness**, cheap**ness**, forgive**ness**, circuitous**ness**, grand**ness**, orderli**ness**, pretentious**ness**, busi**ness**, green**ness**, …*

Phonological Regularities   Morpho-Syntactic Regularities   Semantic Regularities

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | -ness |

```
        N
       / \
      /   \
    Adj   -ness
```

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |

*cool*          *\*coolth*

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |

*tru**th**, warm**th**, wid**th**, dep**th**, fil**th**, slo**th**, streng**th**, dea**th**, dear**th**, weal**th**, leng**th**, you**th**, …*

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |

*truth, warmth, width, depth, filth, sloth, strength, death, dearth, wealth, length, youth,* …

Adj ✗ -th
N

# The Problem of Lexical Uncertainty

| | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

*pine-scented*          *\*pine-scentedity*
*cool*                        *\*coolity*

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

N

Adj ~~X~~ ity

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

*-ile, -al, -able, -ic, -(i)an*

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

*sequentiable*            *sequentiability*

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

N

✗ Adj  ity

N

V  *-ability*  **?**

# The Problem of Lexical Uncertainty

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

# The Problem of Lexical Uncertainty

What principles can the learner use to distinguish real lexical units like *-ness* from false lexical units like *-th*?

| | |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

# Overview

Part 1: An approach to lexical uncertainty.

Part II:  Detailed case study in morphology.

Part III:  Unsupervised models of language learning.

# Overview

Part 1: An approach to lexical uncertainty.

Part II: Detailed case study in morphology.

Part III: Unsupervised models of language learning.

# Bayesian Approach

- Use probabilistic models to define prior distribution over possible lexicons.

- Find posterior distribution over lexicons given actual data using probabilistic conditioning.

$$P(L \mid D) \propto P(D \mid L)\, P(L)$$

# The Framework

1. An underlying computational system that defines space of possible structures.

2. Manage uncertainty over the ways in which forms can be analyzed with a probabilistic model implementing a **tradeoff** between storage and computation (prior and likelihood).

3. Use probabilistic inference to derive language-specific predictions.
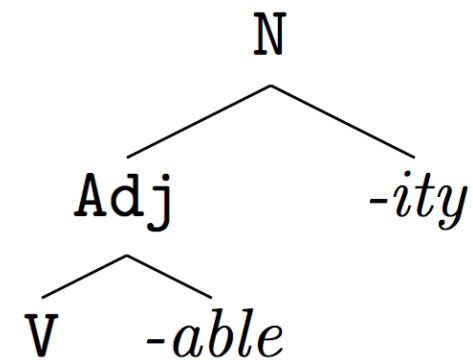
# The Framework

➡️ 1. An underlying computational system that defines space of possible structures.

2. Manage uncertainty over the ways in which forms can be analyzed with a probabilistic model implementing a **tradeoff** between storage and computation (prior and likelihood).

3. Use probabilistic inference to derive language-specific predictions.

# Underlying Computational System

1. Inventory of structured units (lexicon).

2. Structure-building operations.

# Lexical Items

```
         N
        / \
      Adj  -ity
```

Stems, suffixes, prefixes with combinatorial restrictions.

# Structure Building

N

Unit-driven selection.

```
        N
       / \
    Adj   -ity
```

# Structure Building



N

Unit-driven selection.

Adj    *-ity*

V    *-able*

# Structure Building

N    Unit-driven selection.

Adj    -*ity*

V    -*able*

|

*agree*

# The Framework

1. An underlying computational system that defines space of possible structures.

2. Manage uncertainty over the ways in which forms can be analyzed with a probabilistic model implementing a **tradeoff** between storage and computation (prior and likelihood).

3. Use probabilistic inference to derive language-specific predictions.

# Uncertainty

Phonological, semantic, and morphosyntactic processes generate candidate analyses.

N
Adj            *-ity*
V      *-able*
|
*agree*

# Uncertainty



```
        N
      /   \
    Adj    -ity
   /   \
  V     -able
  |
agree
```

# Uncertainty



N

|

agreeability

# Uncertainty

# Probabilistic Model

- **Goal**: Define a probability distribution that quantifies uncertainty about likely and unlikely lexical units.

- Prior and likelihood lead to a **tradeoff** between storage and computation.

  1. Prior over lexicon: P(L).

  2. Likelihood over derivations of individual forms: P(D | L).

# Prior on Lexicon

- Prefer a small number of highly reusable (i.e., generalizable) units.

- Dirichlet Process (Ferguson, 1973; Aldous, 1983; Sethuraman, 1994; Pitman 1995).

  - Allows unbounded number of stored units (non-parametric).

  - Prefers fewer units.

  - Prefers more reusable units.

$$P(\text{unit}) \propto \text{frequency of use}$$

# Likelihood of Derived Forms

- Simple observed forms with few parts.

- Probability of derivation: Product of probability of individual units.

  - Probabilities between $0$ and $1$.

    - Geometric decrease in probability.

  - Prefers fewer lexical items per derivation.

  - Prefers to store more complex units.

# The Framework

1. An underlying computational system that defines space of possible structures.

2. Manage uncertainty over the ways in which forms can be analyzed with a probabilistic implementing a **tradeoff** between storage and computation.

3. Use probabilistic inference to derive language-specific predictions.

# The Mathematical Model: *Fragment Grammars*

- Generalization of *Adaptor Grammars* (Johnson et al., 2007).

  - Allows storing of partial trees.

- Can be generalized with stochastic variants of *memoization* and *order of evaluation*.

# Inference Problem

Search for sets of stored units and derivations of individual forms which best explain the input taking into account preferences for a small inventory of highly reusable units and simple derivations of individual forms.

# Example Input

# Maximal Decomposition

# Derivation Complexity

# Derivation Complexity

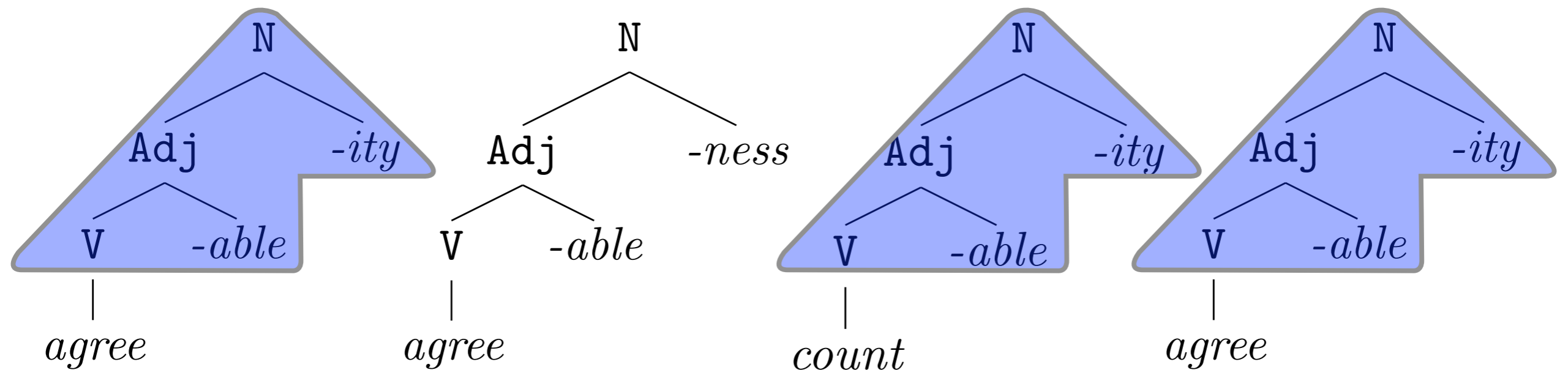$$P\left( \begin{array}{c} \text{N} \\ \text{Adj} \quad \text{-ity} \end{array} \right)$$

N

Adj          -ity

V          -able

agree

# Derivation Complexity

$$P(\text{[Adj -ity → N]}) \times P(\text{[V -able → Adj]})$$

# Derivation Complexity

# Sharing Across Expressions

# Sharing Across Expressions

# No Decomposition

# Derivation Complexity

```
                    N
                   / \
                Adj   -ity
               /  \
              V   -able
              |
            agree
```

# Derivation Complexity

# Sharing Across Expressions

# Sharing Across Expressions

# Intermediate Decomposition

# Derivation Complexity

# Derivation Complexity



P(
N
Adj        -ity
V    -able
)

N
Adj        -ity
V    -able
|
*agree*

# Derivation Complexity

# Sharing Across Expressions

# Sharing Across Expressions

# Computation/Storage Tradeoff

# Computation/Storage Tradeoff



Builds on both of classic and recent work in non-parametric Bayesian and minimum description length framework.s

Has connections to theory of programming languages, Kolmogorov complexity, and other areas, via probabilistic programming and program induction frameworks.

# Overview

Part 1: An approach to lexical uncertainty.

Part II:  Detailed case study in morphology.

Part III:  Unsupervised models of language learning.

# English Derivational Morphology

|  | Suffix |
|---|---|
| Productive | -ness |
| Unproductive | -th |
| Productive in Combination | -ity |

# Morphological Productivity

1. **Background**

2. **Productivity and frequency**

3. **Affix ordering**

# Morphological Productivity

1. **Background**

2. **Productivity and frequency**

3. **Affix ordering**

# Idealized Learner Study

- Not intended to correspond to a particular stage of development or learning procedure used by children.

- **Question**: Does the computation/storage tradeoff provide an effective way of distinguishing between productive units like *-n*ess and false generalizations like *-th*?

# Input Representations

Phonologically, semantically, and morphosyntactically plausible candidate analyses, erring on the side of shareable structure (decomposition).

# Models



*Full-Parsing*
(DMPCFG)

*Full-Storage*
(MAG)

*Inference-Based*
(FG)

1. Differ in storage strategy, otherwise as closely matched as possible.

2. Same inputs and representation space.

3. Represent state-of-the-art implementations from natural language processing.

# Models



Full-Parsing (DMPCFG)

Full-Storage (MAG)

Inference-Based (FG)

# *Full-Parsing*

(MAP Dirichlet-Multinomial Context-Free Grammars: `DMPCFG`)

Johnson, Griffiths, & Goldwater (2007a)



- <u>Storage Strategy</u>: Maximally Decompose.

- Implements radical version of decompositional theories from psychology (e.g., Taft, 1988) and linguistics (e.g., Halle & Marantz, 1983).

# Full-Parsing

## (MAP Dirichlet-Multinomial Context-Free Grammars: `DMPCFG`)

Johnson, Griffiths, & Goldwater (2007a)



- **<u>Storage Strategy</u>**: Maximally Decompose.

- Best thought of as a baseline system for morphology.

# *Full-Parsing*

(MAP Dirichlet-Multinomial Context-Free Grammars: `DMPCFG`)

Johnson, Griffiths, & Goldwater (2007a)

- <u>Storage Strategy</u>: Maximally Decompose.

Productivity

$$P\left(\begin{smallmatrix}N\\Adj \quad \text{-}ness\end{smallmatrix}\right) \propto \text{Token Frequency}$$

# Models



**Full-Parsing** (DMPCFG)

**Full-Storage** (MAG)

**Inference-Based** (FG)

# Full-Storage

## (MAP All-Adapted Adaptor Grammars: MAG)

Johnson, Griffiths, & Goldwater (2007)



- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.
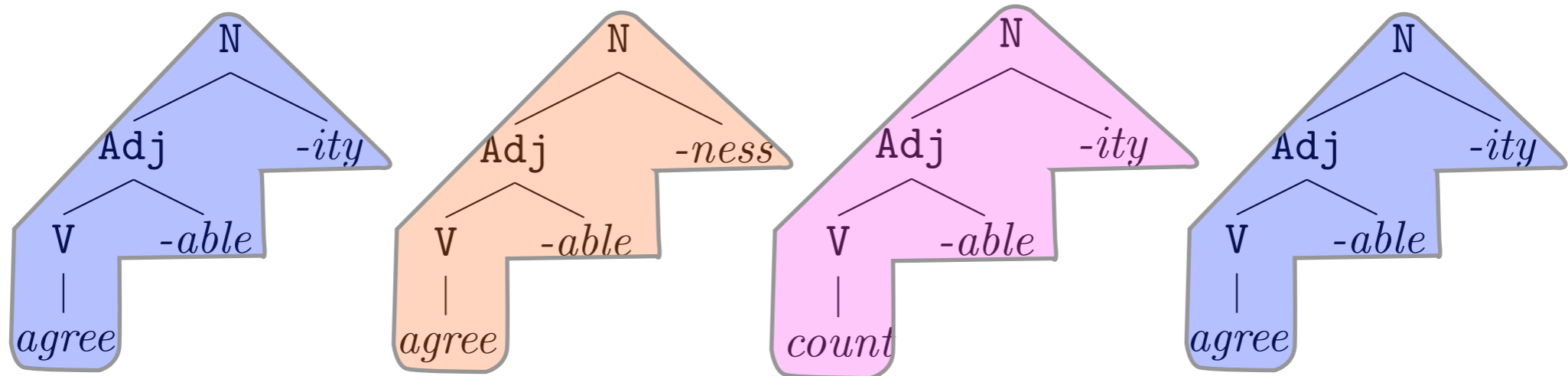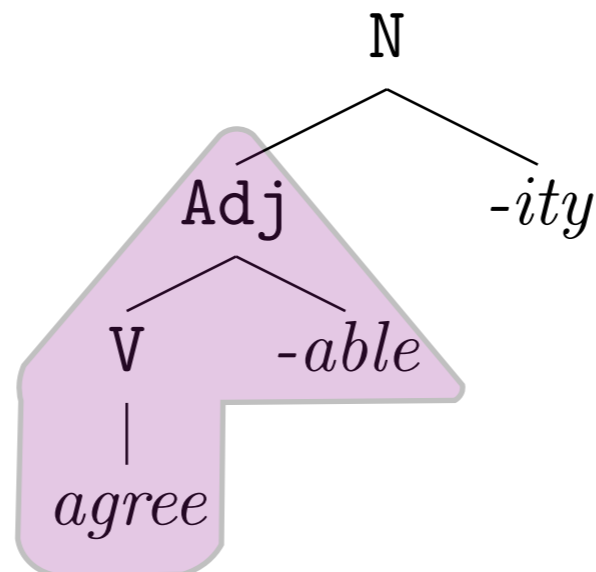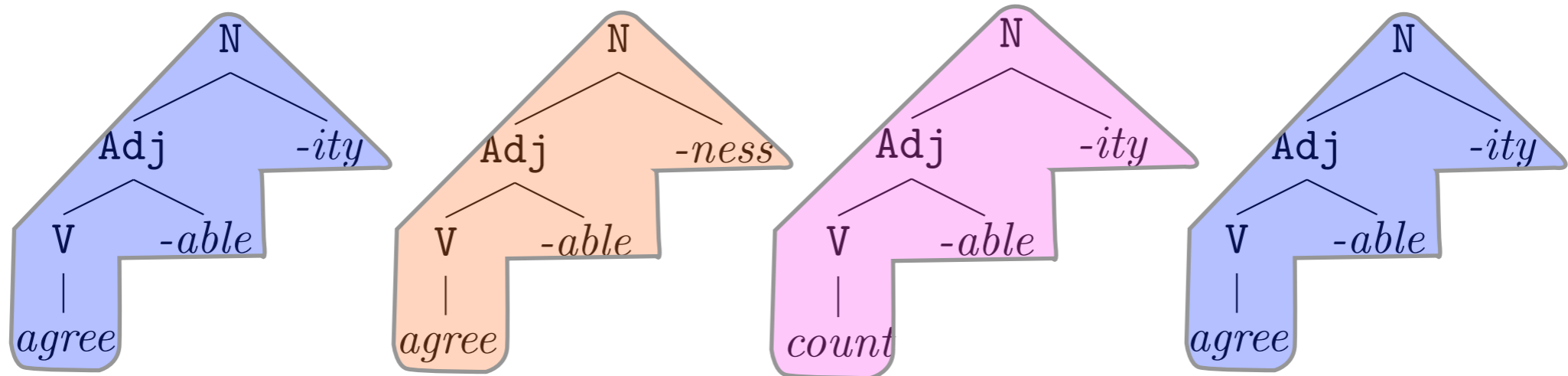
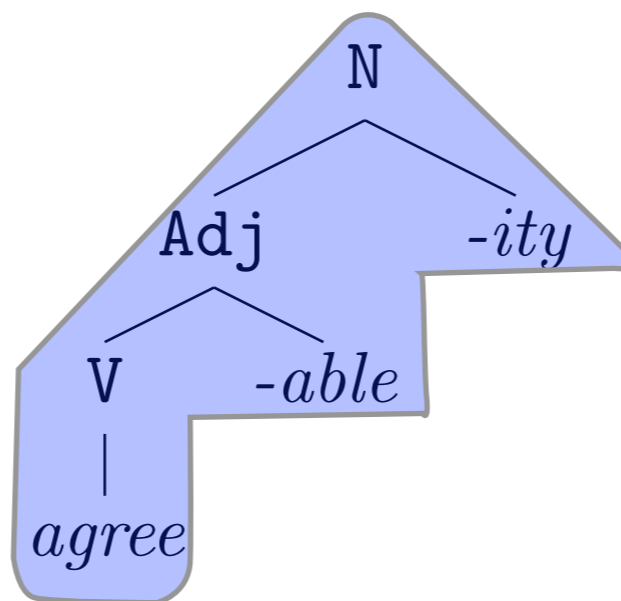# *Full-Storage*

## (MAP All-Adapted Adaptor Grammars: MAG)

Johnson, Griffiths, & Goldwater (2007)



- **<u>Storage Strategy</u>**: Store structures in their entirety after first time generated.

All full sub-trees are stored recursively.

# *Full-Storage*

(MAP All-Adapted Adaptor Grammars: `MAG`)

Johnson, Griffiths, & Goldwater (2007)



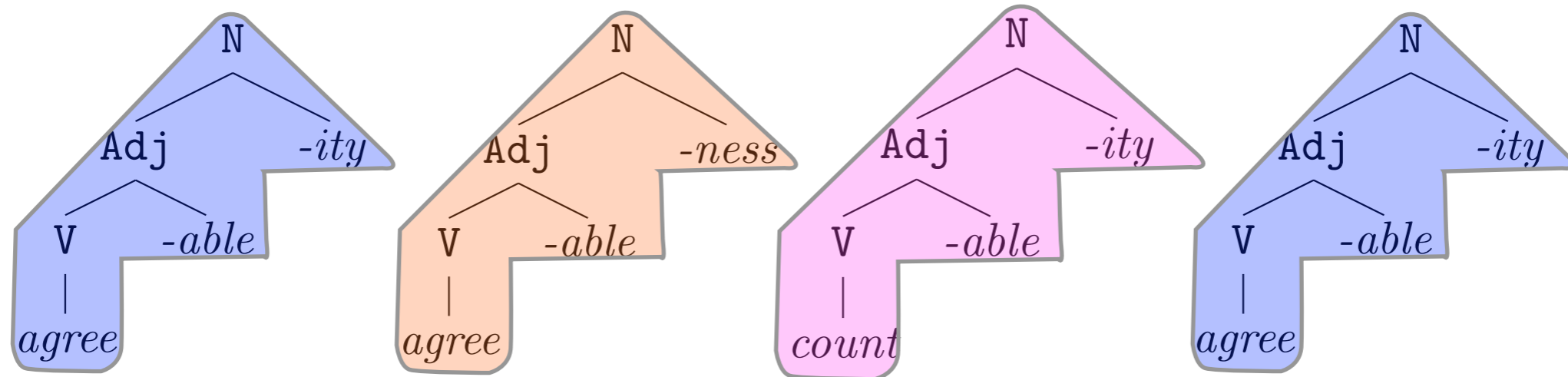- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.

All full sub-trees are stored recursively.

# *Full-Storage*

(MAP All-Adapted Adaptor Grammars: `MAG`)
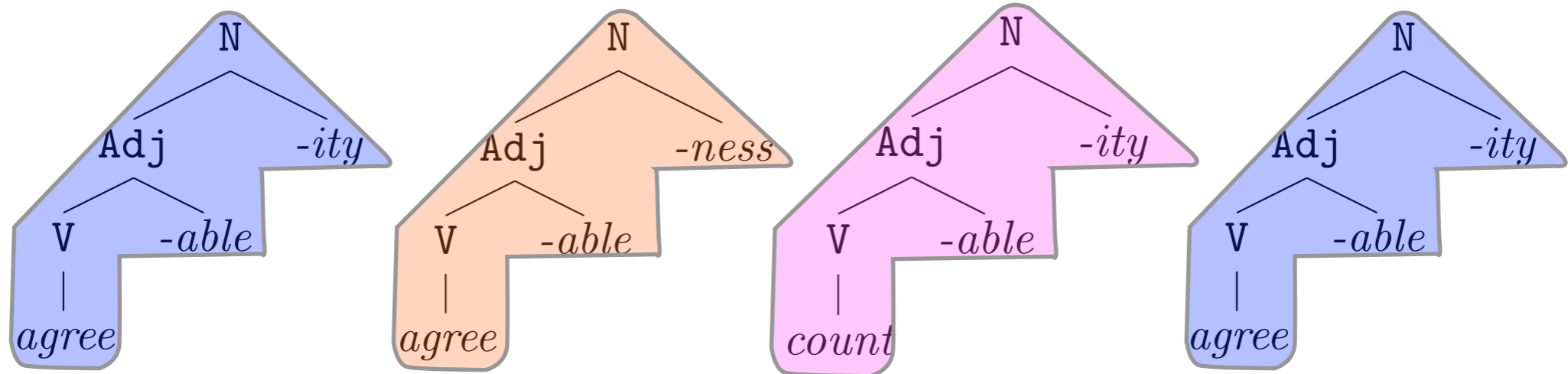
Johnson, Griffiths, & Goldwater (2007)



- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.

All full sub-trees are stored recursively.

# Full-Storage

(MAP All-Adapted Adaptor Grammars: MAG)

Johnson, Griffiths, & Goldwater (2007)



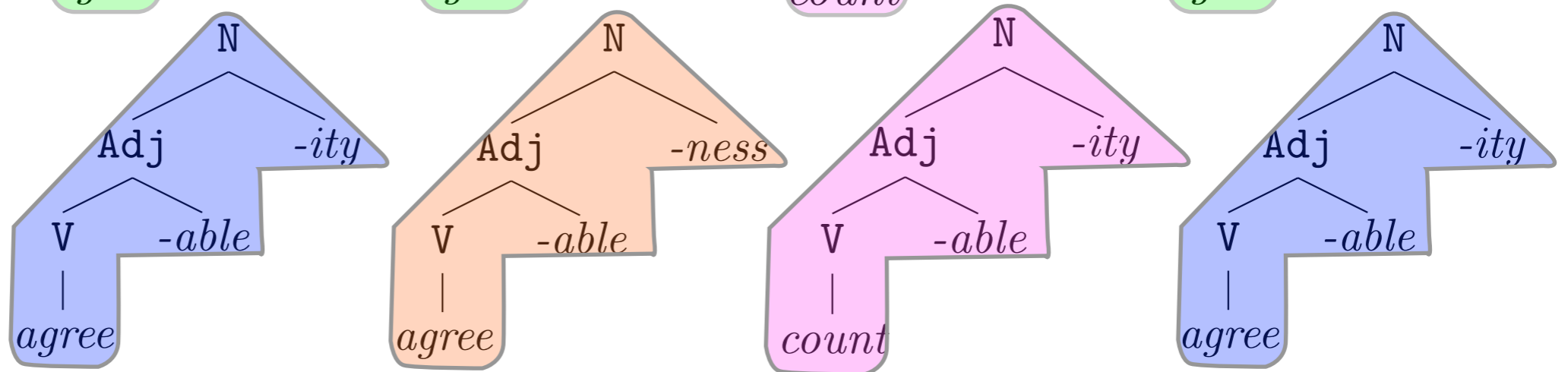- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.

All full sub-trees are stored recursively.

# *Full-Storage*

## (MAP All-Adapted Adaptor Grammars: `MAG`)

Johnson, Griffiths, & Goldwater (2007)



- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.

- Modern probabilistic implementation of classical lexical redundancy rules (e.g., Jackendoff, 1975; Aronoff, 1976).

# Full-Storage

## (MAP All-Adapted Adaptor Grammars: MAG)

Johnson, Griffiths, & Goldwater (2007)



- <u>Storage Strategy</u>: Store structures in their entirety after first time generated.

## Productivity

$$P\left(\; \underset{\text{Adj} \quad \text{-ness}}{\overset{\text{N}}{\triangle}} \;\right) \;\propto\; \text{Type frequency}$$

# Models



**Full-Parsing** (DMPCFG)

**Full-Storage** (MAG)

**Inference-Based** (FG)

# *Inference-Based*

## (Fragment Grammars: FG)

O'Donnell et al. (2009); O'Donnell (2011, 2015)



- <u>Storage Strategy</u>: Store set of units that best explains data (only inference-based storage proposal).
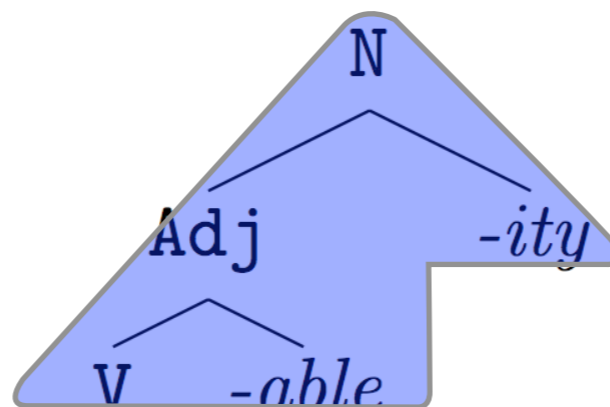
# *Inference-Based*

## (Fragment Grammars: FG)
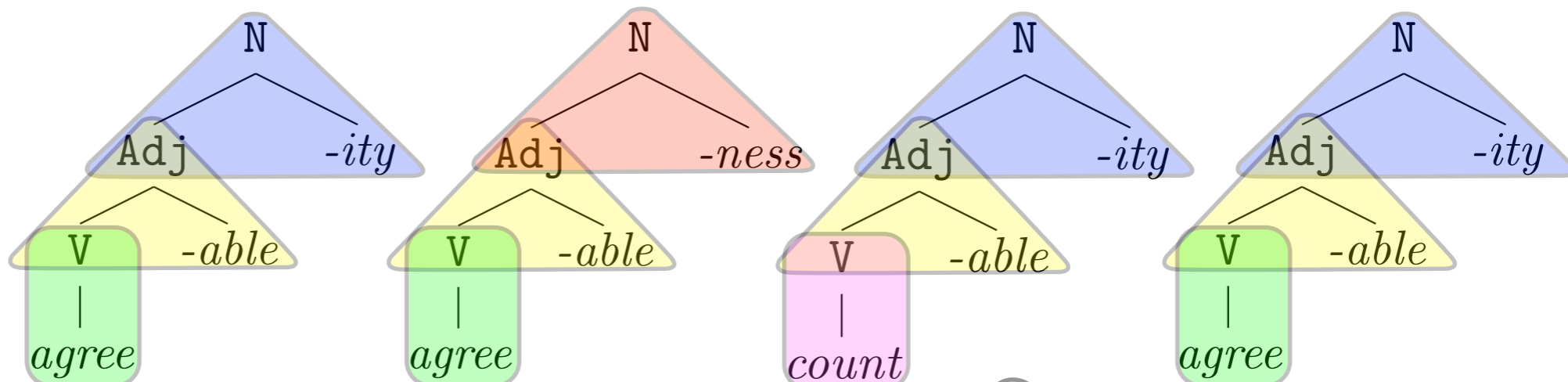O'Donnell et al. (2009); O'Donnell (2011, 2015)



- <u>Storage Strategy</u>: Store set of units that best explains data (only inference-based storage proposal).
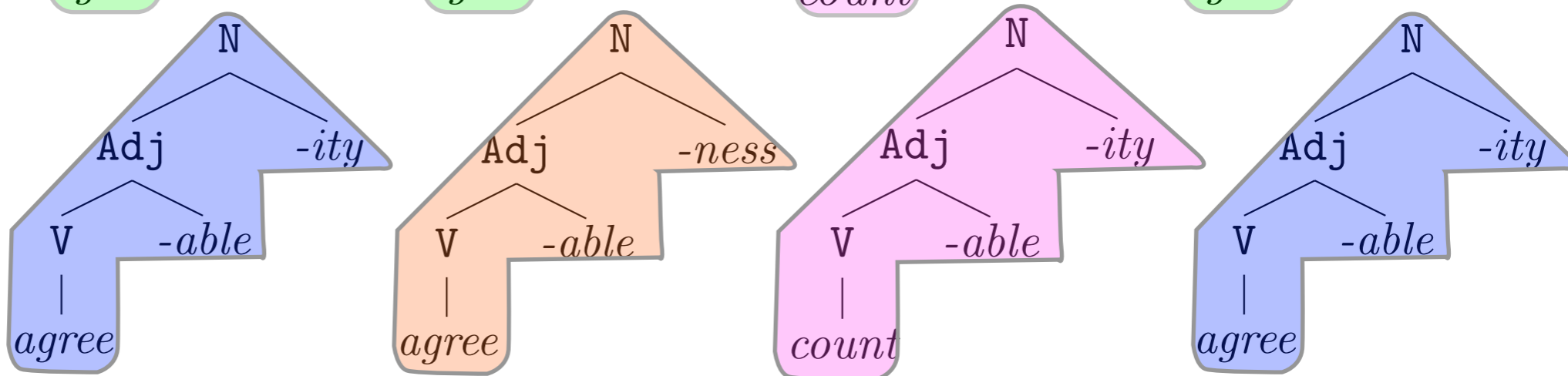
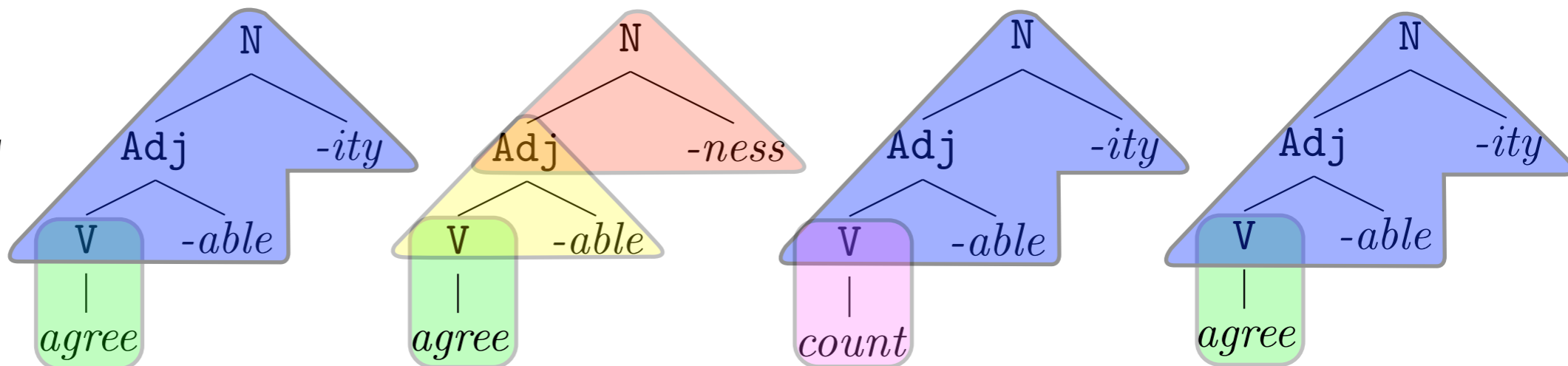- Generalization of *Full-Storage* (AG) model (Johnson et al., 2007).

# Models

# Morphological Productivity

1. **Background**
2. **Productivity and frequency**
3. **Affix ordering**

# Productivity and the Distribution of Forms

What aspects of the distribution of words and part-words signal that some structure is an independent, productive lexical unit?

# Five Most Productive Suffixes (Predicted)

## Full-Parsing (MDPCFG)

| Suffix | Example |
|---|---|
| *ion*:V>N | *regression* |
| *ly*:Adj>Adv | *quickly* |
| *ate*:BND>V | *segregate* |
| *ment*:V>N | *development* |
| *er*:V>N | *talker* |

## Inference-Based (FG)

| Suffix | Example |
|---|---|
| *ly*:Adj>Adv | *quickly* |
| *er*:V>N | *talker* |
| *ness*:Adj>N | *tallness* |
| *y*:N>Adj | *mousey* |
| *er*:N>N | *prisoner* |

## Full-Storage (MAG)

| Suffix | Example |
|---|---|
| *ly*:Adj>Adv | *quickly* |
| *ion*:V>N | *regression* |
| *er*:V>N | *talker* |
| *ly*:V>Adv | *bitingly* |
| *y*:N>Adj | *mousey* |

# Five Most Productive Suffixes (Predicted)

## -ness

### Full-Parsing (MDPCFG)

| Suffix | Example |
|---|---|
| ion:V>N | regression |
| ly:Adj>Adv | quickly |
| ate:BND>V | segregate |
| ment:V>N | development |
| er:V>N | talker |

### Inference-Based (FG)

| Suffix | Example |
|---|---|
| ly:Adj>Adv | quickly |
| er:V>N | talker |
| ness:Adj>N | tallness |
| y:N>Adj | mousey |
| er:N>N | prisoner |

### Full-Storage (MAG)

| Suffix | Example |
|---|---|
| ly:Adj>Adv | quickly |
| ion:V>N | regression |
| er:V>N | talker |
| ly:V>Adv | bitingly |
| y:N>Adj | mousey |

# Five Most Productive Suffixes
# (Predicted)

## *-ion*

| *Full-Parsing* (MDPCFG) | *Inference-Based* (FG) | *Full-Storage* (MAG) |

| Suffix | Example |
|--------|---------|
| *ion*:V>N | *regression* |
| *ly*:Adj>Adv | *quickly* |
| *ate*:BND>V | *segregate* |
| *ment*:V>N | *development* |
| *er*:V>N | *talker* |

| Suffix | Example |
|--------|---------|
| *ly*:Adj>Adv | *quickly* |
| *er*:V>N | *talker* |
| *ness*:Adj>N | *tallness* |
| *y*:N>Adj | *mousey* |
| *er*:N>N | *prisoner* |

| Suffix | Example |
|--------|---------|
| *ly*:Adj>Adv | *quickly* |
| *ion*:V>N | *regression* |
| *er*:V>N | *talker* |
| *ly*:V>Adv | *bitingly* |
| *y*:N>Adj | *mousey* |

## *eat : *eation*

# Productivity and Frequency Distributions

- Productive suffixes give rise to new forms.

- <u>Distributional consequence</u>: Large proportion of low-frequency forms.

- Large number of rare event distributions (LNRE) investigated mathematically (Khmaladze, 1987; Baayen, 2001).

- Used to develop useful statistical tools to quantify and study productivity in corpus samples (e.g., Baayen, 1992).

  - More below.

# Productivity and The Proportion of Low-Frequency Forms

*-ness*

*-ion*



34% Freq-One Forms

.05% Freq-One Forms

Number of Words

Word Frequency (Log Scale)

Type and Token Frequencies

# Five Most Productive Suffixes (Predicted)

*-ne*

High Proportion of Low-Frequency Forms

*Full-Parsing* (MDPCFG)

| Suffix | Example |
|---|---|
| *ion*:V>N | *regression* |
| *ly*:Adj>Adv | *quickly* |
| *ate*:BND>V | *segregate* |
| *ment*:V>N | *development* |
| *er*:V>N | *talker* |

*Inference-Bas...* e (MAG)

| Suffix | Example |
|---|---|
| *ly*:Adj>Adv | *quickly* |
| *er*:V>N | *talker* |
| *ness*:Adj>N | *tallness* |
| *y*:N>Adj | *mousey* |
| *er*:N>N | *prisoner* |

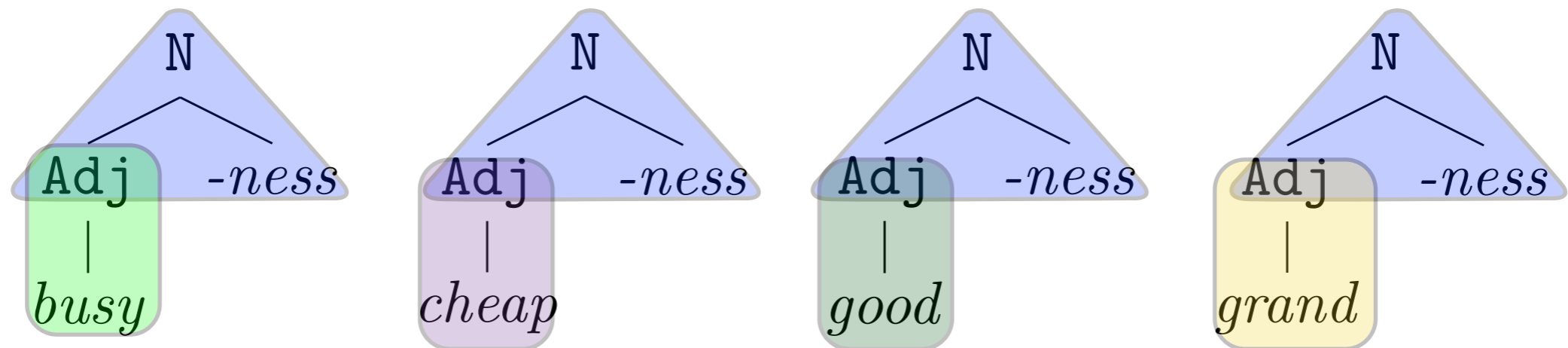| Suffix | Example |
|---|---|
| *ly*:Adj>Adv | *quickly* |
| *ion*:V>N | *regression* |
| *er*:V>N | *talker* |
| *ly*:V>Adv | *bitingly* |
| *y*:N>Adj | *mousey* |

# Productivity and The Proportion of Low-Frequency Forms

- The *Inference-Based* (FG) makes correct use of distributional facts about productivity.
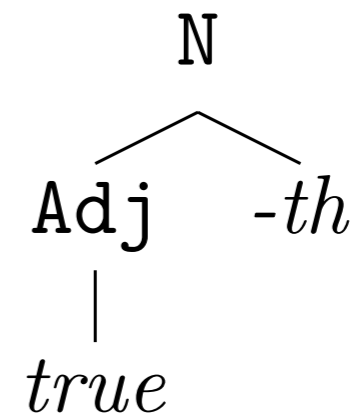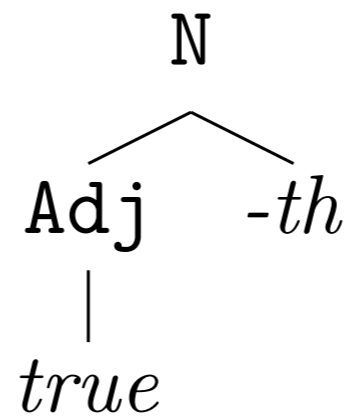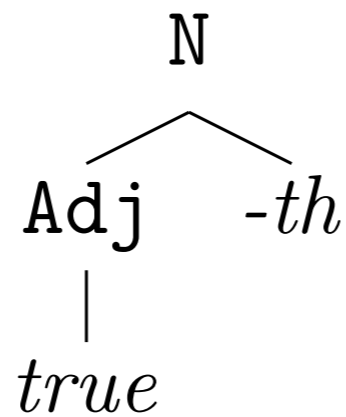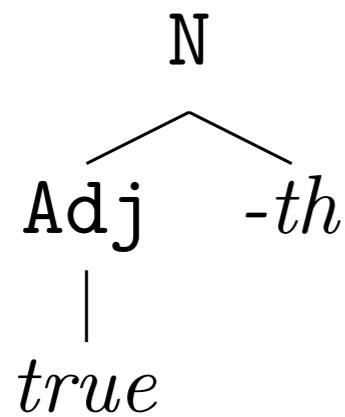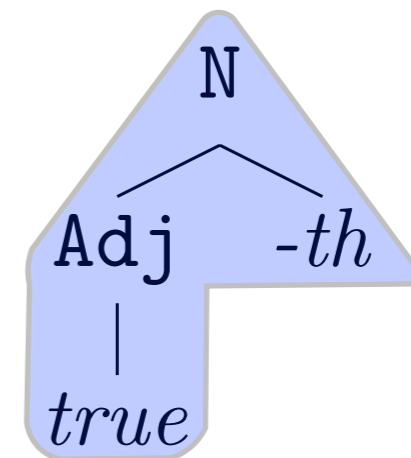
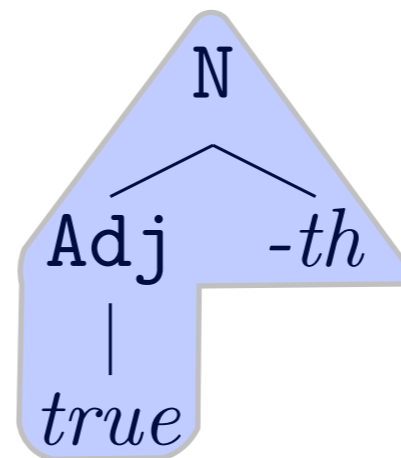- Why is the model sensitive to the proportion of low-frequency forms?
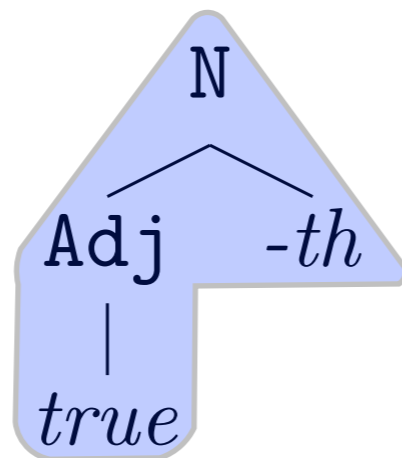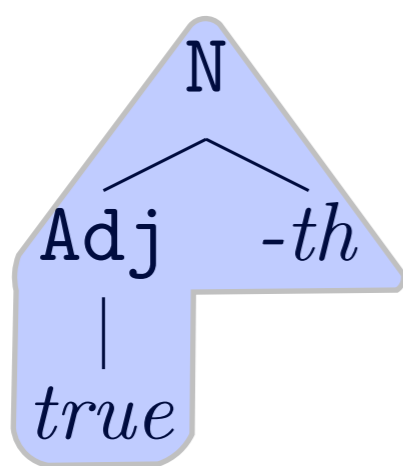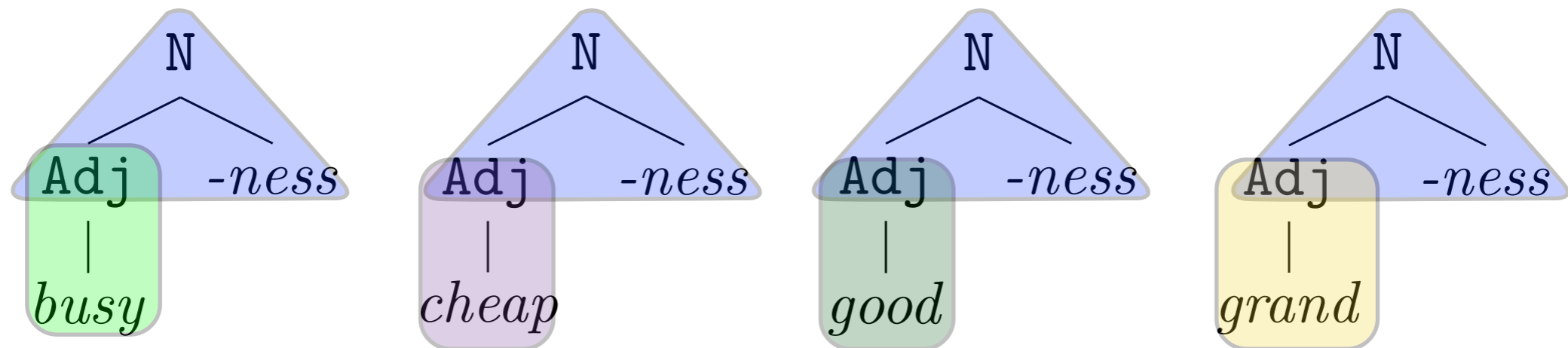
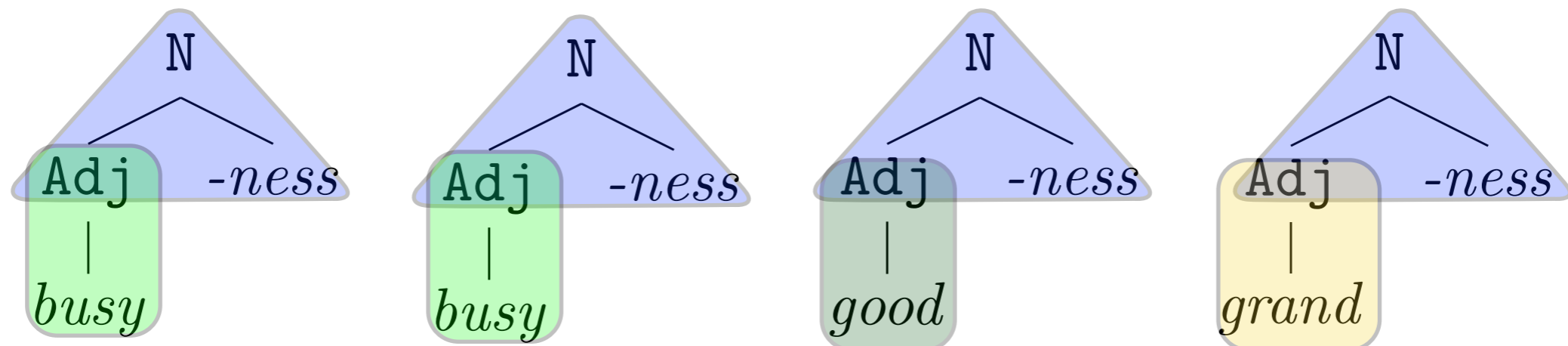# Ideal Input for Productivity

# Ideal Input for Productivity

# Ideal Input for Storage

# Ideal Input for Storage

# Productivity and the Proportion of Low-Frequency Forms

# Productivity and the Proportion of Low-Frequency Forms

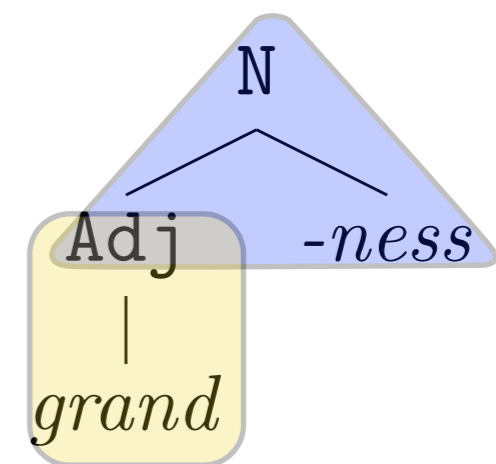# Productivity and the Proportion of Low-Frequency Forms
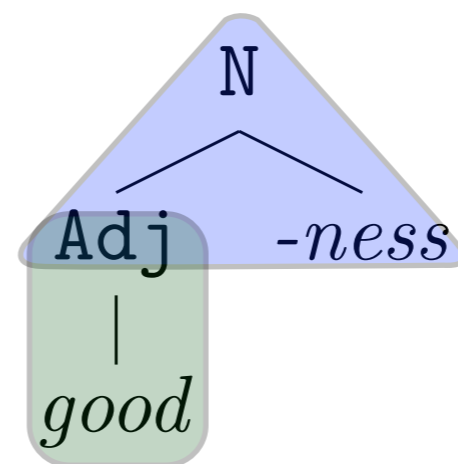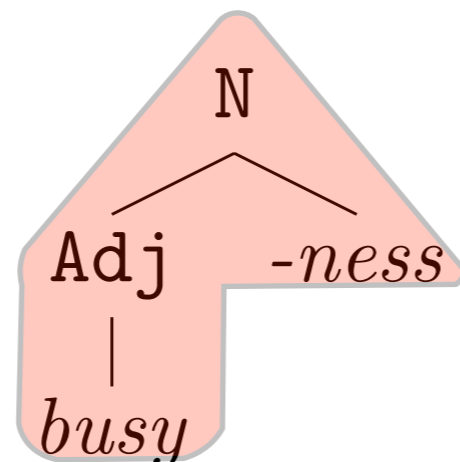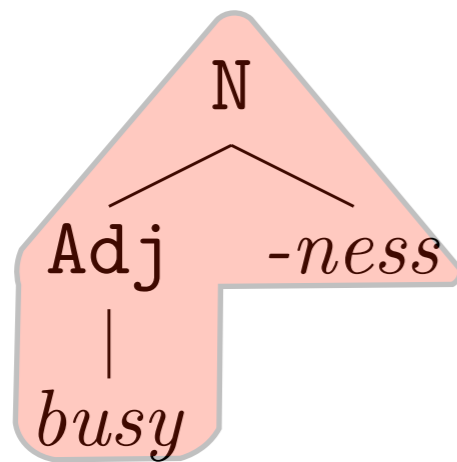
# Productivity and the Proportion of Low-Frequency Forms

# Productivity and the Proportion of Low-Frequency Forms

# Baayen's Hapax-Based Estimators

- Baayen's $\mathcal{P} / \mathcal{P}^*$ (e.g., Baayen, 1992)

- Most well-studied statistical corpus-oriented estimators of productivity.

- Estimators of *single-affix* productivity based on the proportion hapaxes.

- Various mathematical derivations.

  - Rate of vocabulary change.

  - Good-Turing.

- Non-psychological.

# The *Inference-Based* Model and Baayen's Measures

- For individual suffixes predictions strongly correlate.

| Measure | FG | MDPCFG | MAG |
|---------|-----|--------|-----|
| $\mathcal{P}$ | **0.907** | -0.0003 | 0.692 |
| $\mathcal{P}*$ | **0.662** | 0.480 | 0.568 |

- *Inference-Based* (FG) model makes no explicit assumption about-low frequency forms.

- *Inference-Based* (FG) model derives relationship between productivity inferences and distributional facts from storage-computation tradeoff applied to the problem of lexicon learning.

# Morphological Productivity

# The Affix-Ordering Problem



- Tiny fraction of logically possible affix combinations occur in practice.

- Fabb (1988)
  - 43 suffixes
  - 663 predicted possible.
  - ~50 attested.

# The Affix-Ordering Problem

- **Numerous accounts** (see, O'Donnell, 2015 for review).

  - None entirely successful.

- <u>Main empirical generalization</u>: Ordering correlates with a number of other properties.

# Correlated Properties

Earlier Suffixes

Later Suffixes

*-ion, -ity, -y, -al, -ic,*
*-ate, -ous, -ive, …*

*-ness, -less, -hood, -ful,*
*-ly, -y, -like, -ist, …*

1. Origin (Latin/Germanic)
2. Phonological regularity
3. Transparency of meaning
4. Productivity

# Productivity and Ordering Generalization

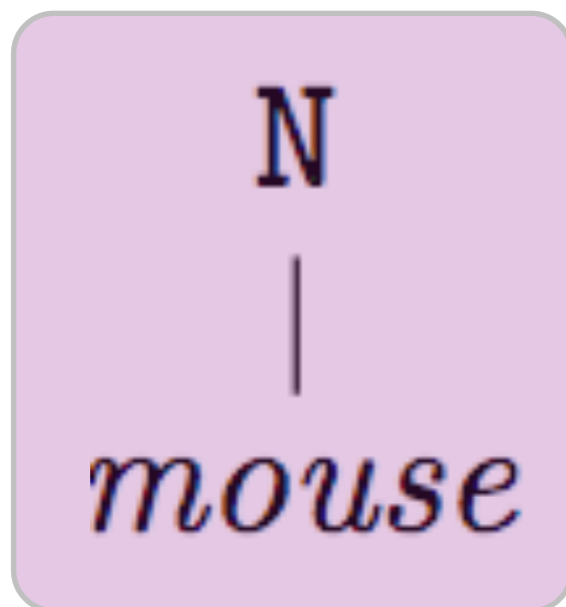On average more productive affixes appear outside of less productive affixes
(Hay, 2002; Hay & Plag, 2004; Plag et al., 2009).
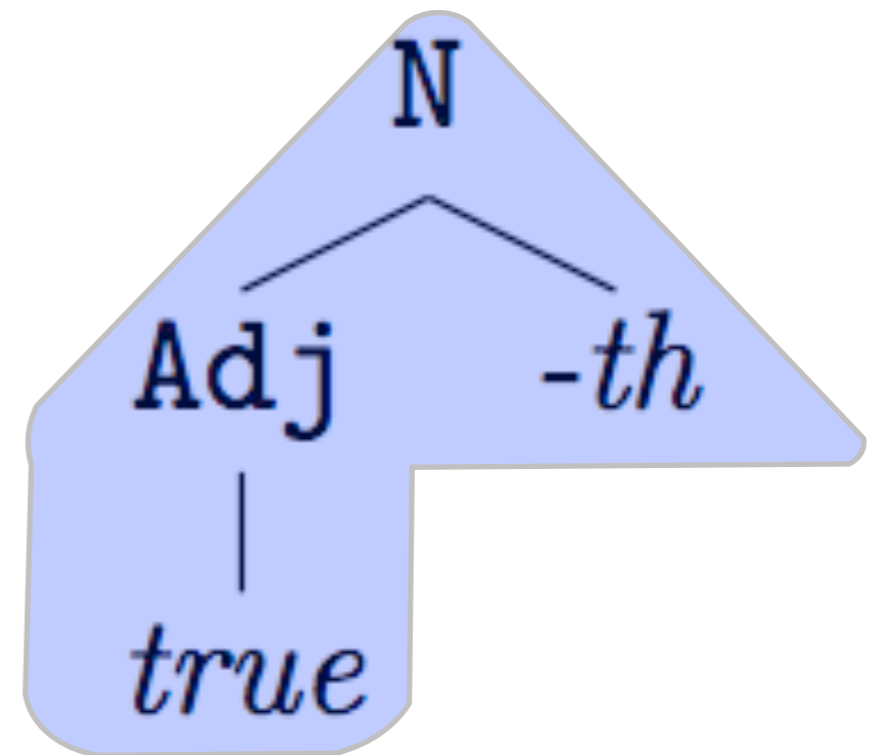
# Productivity and Ordering Generalization

- Follows as a consequence of *Inference-Based* (FG) model's pattern of storage and computation.

- **Crucial point**: By definition, productive affixes are just independent units which can combine freely; unproductive affixes are subparts of other stored forms.
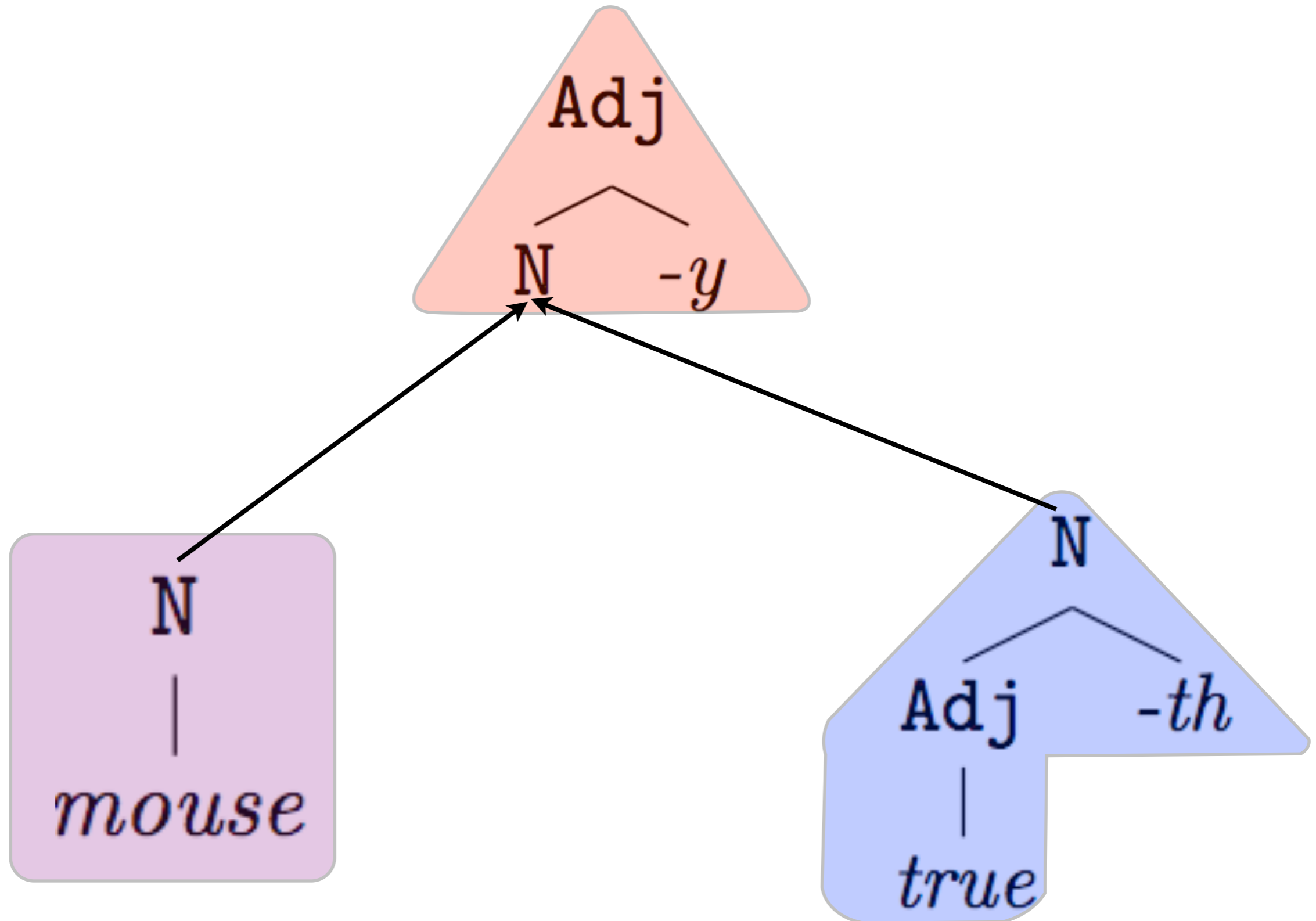
# Productivity and Ordering

Simple words

Words with
unproductive suffixes

# Productivity and Ordering

# Correlated Properties

- <u>Latinate v. Germanic</u> reflects diachronic fact that Latinate affixes were borrowed as parts of loan words.

- <u>Regularity and Transparency</u> differences are accounted for if stored items can preserve or accrue idiosyncrasies.

# Paradoxical Suffix Combinations

- Exceptions to the productivity and ordering generalization and all current theoretical accounts (Hay, 2003).

    - *-ability*, *-ation*, *-istic*

- **Idea**: *-ability*, *-ation*, and *-istic* are single, stored lexical units.

# Predicted Generalizable Units

| Sequence | Category |
|----------|----------|
| -ate -ion | N |
| -ic -al | A |
| -ate -ive | A |
| -al -ity | N |
| -al -ize | V |
| -ology -ist | N |
| -ment -al | A |
| -able -ity | N |
| -ist -ic | A |
| -ous -ity | N |

# Experimental Evidence

(Aronoff & Schvaneveldt,1978; Anshen and Aronoff, 1981)

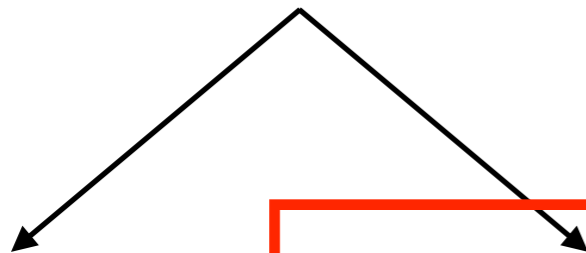| Suffix 1 | Suffix 2 |
|---|---|
| *-ive   v.  -(a)ble* | *-ity   v.  -ness* |

*novel stem* +
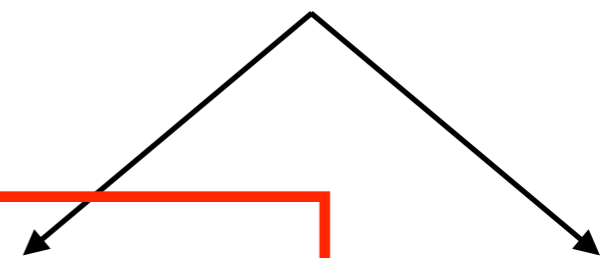
*depulsive*
(Aronoff & Schvaneveldt, 1978)

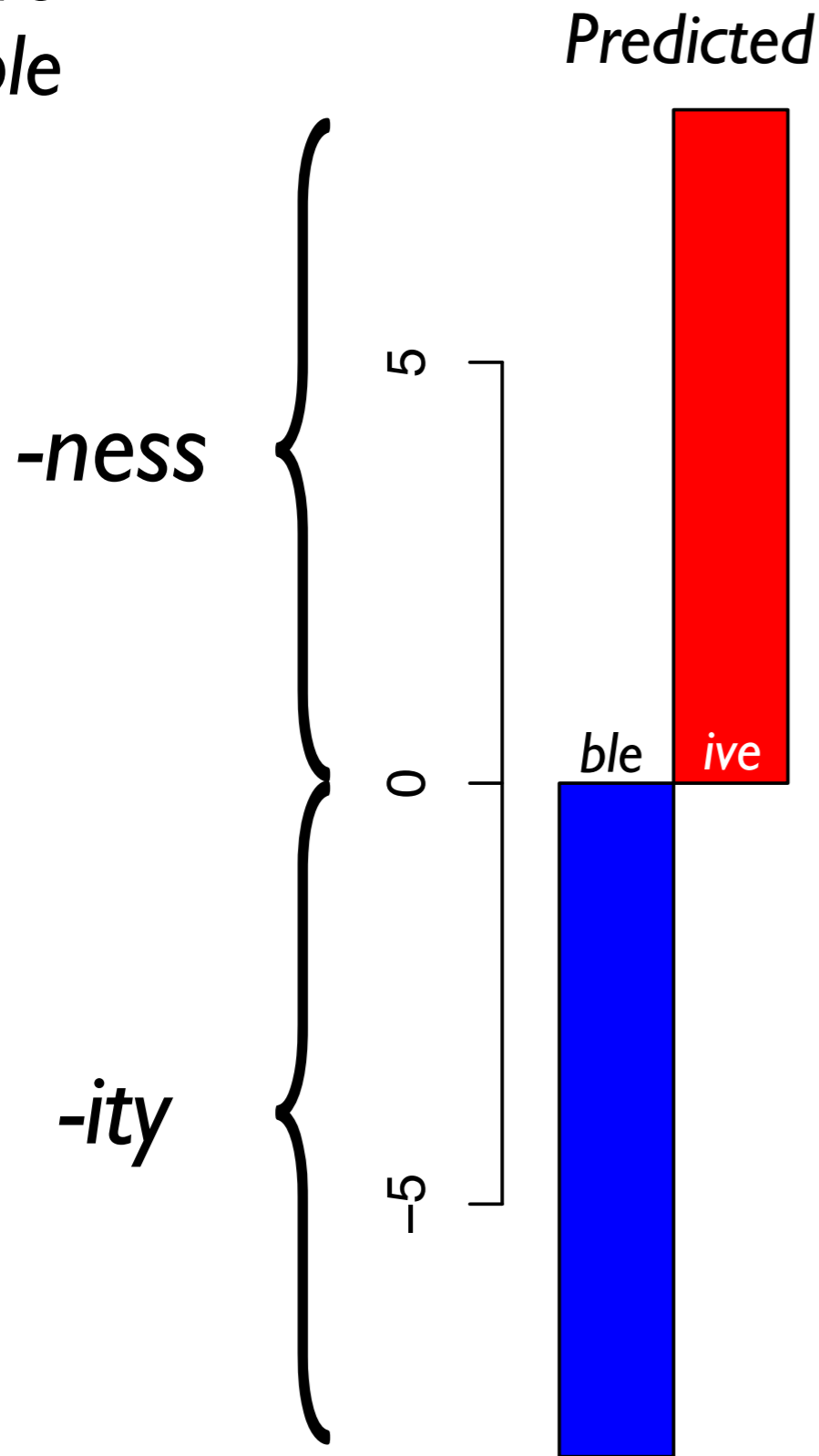*depulsivity*   depulsiveness

*remortible*
(Anshen & Aronoff, 1981)

remortibility   *remortibleness*
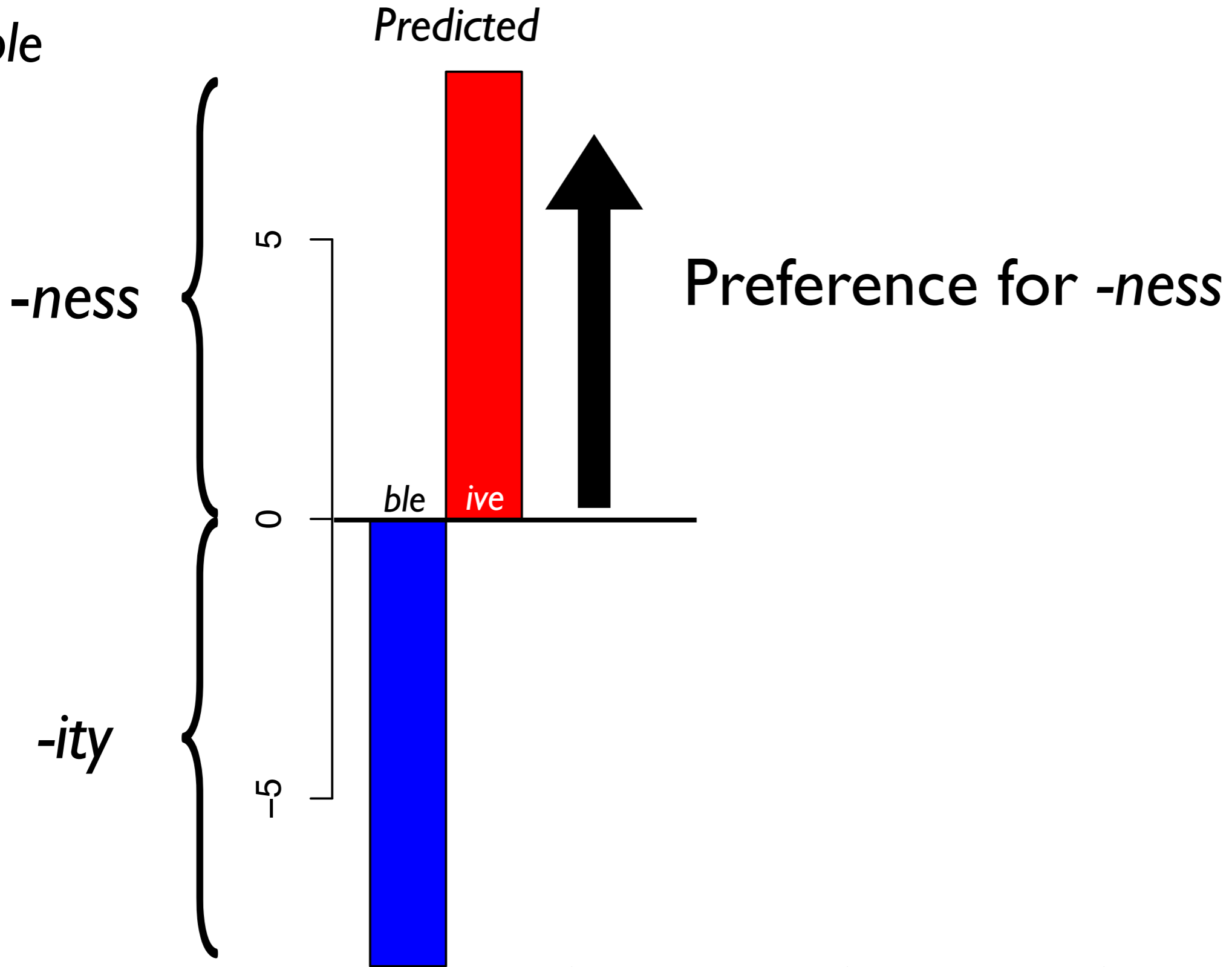
# -ivity v. -bility



-ive
-ble

-ness

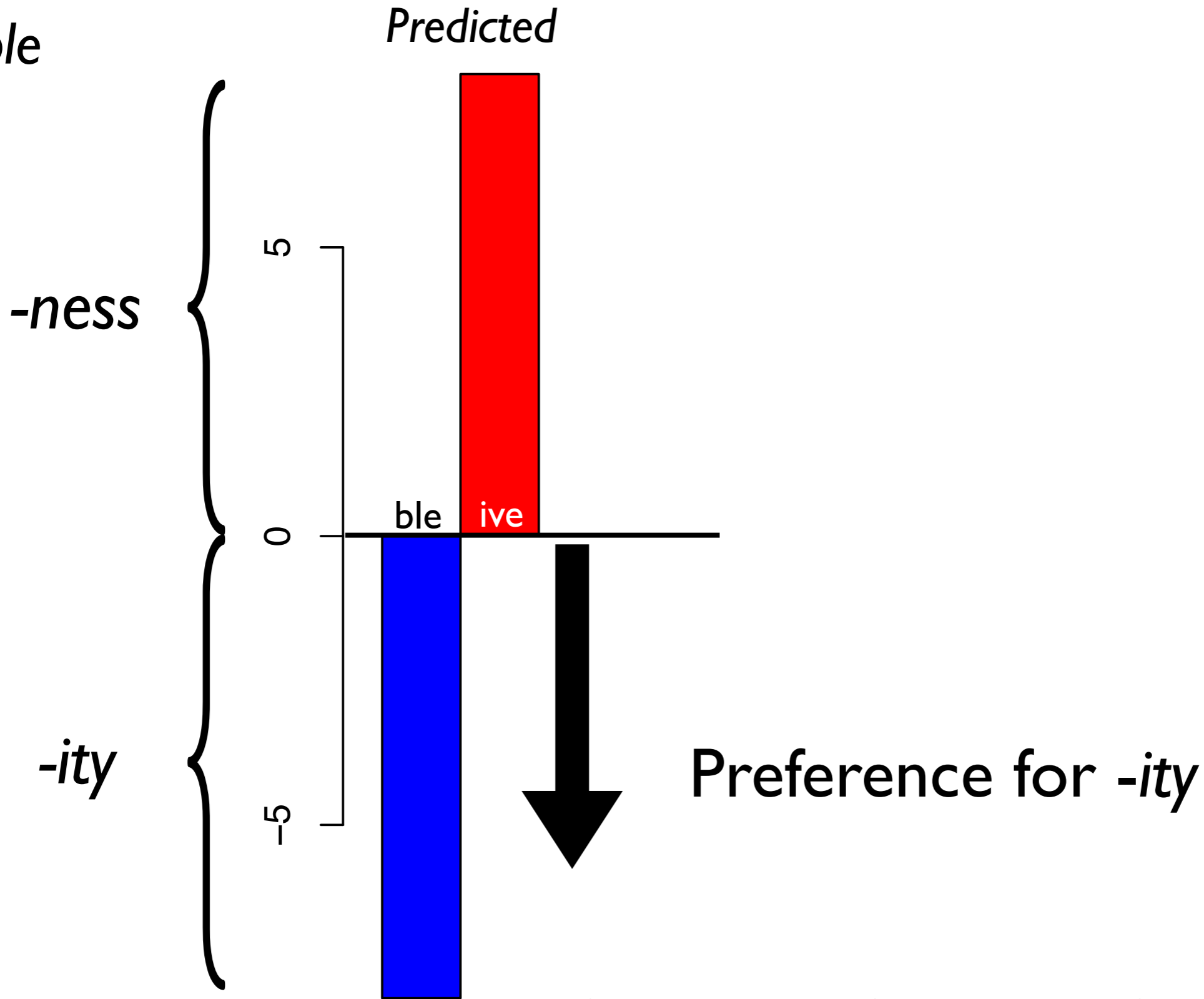-ity

Predicted

Preference for -ness

ble    ive

5

0

−5

# *-ivity* v. *-bility*

# Inference-Based (FG)

# Overview

Part 1: An approach to lexical uncertainty.

Part II:  Detailed case study in morphology.

Part III:  Unsupervised models of language learning.

# Morphology Studies

Phonological, morphosyntactic, and semantic processes generate candidate analyses.

N

Adj          -*ity*

V     -*able*

|

*agree*

# Morphology Studies

Phonological, morphosyntactic, and semantic processes generate candidate analyses.

```
            N
           / \
        Adj   -ity
       /   \
      V    -able
      |
    agree
```

# Unsupervised Learning

<span style="color:red">Phonological</span>, <span style="color:red">morphosyntactic</span>, and semantic processes generate candidate analyses.

N
 ├── Adj
 │     ├── V
 │     │    └── *agree*
 │     └── *-able*
 └── *-ity*

# Unsupervised Learning

Phonological, morphosyntactic, and semantic processes generate candidate analyses.

No Productivity Inference



*agree*     *-able*     *-ity*

# Unsupervised Lexicon Discovery from Speech
### (Lee, O'Donnell, & Glass, 2015)

- First step in integrating models of phonetic, phonological, and morphological structure learning.

- Completely unsupervised learning of words, and morphemes from speech (acoustic input).

- Uses similar storage-computation tradeoffs in each component.

# Unsupervised Lexicon Discovery from Speech
(Lee, O'Donnell, & Glass, 2015)

1. Model of unsupervised learning of phonological units from acoustic data.

2. Model of mapping between underlying (lexical) to surface phonological units.

3. Model of morphological structure.
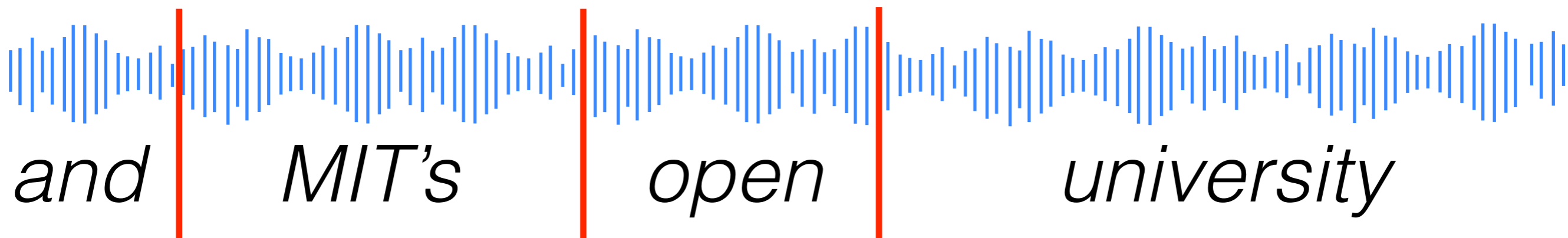
4. Model of lexical storage.

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech

(Lee, O'Donnell, & Glass, 2015)



*and* | *MIT's* | *open* | *university*

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

> 1. **Learning Phonemes** (Lee & Glass, 2012)
>
>    - Phones are hidden Markov models generating acoustics.
>    - Unbounded number of phones.
>    - Dirichlet process prior over phone inventory.
>      - Fewer more reusable phones.

*and* | *MIT's* | *open* | *university*

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)



**1.Learning Phonemes** (Lee & Glass, 2012)

/b/

*and* | *MIT's* | *open* | *university*

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
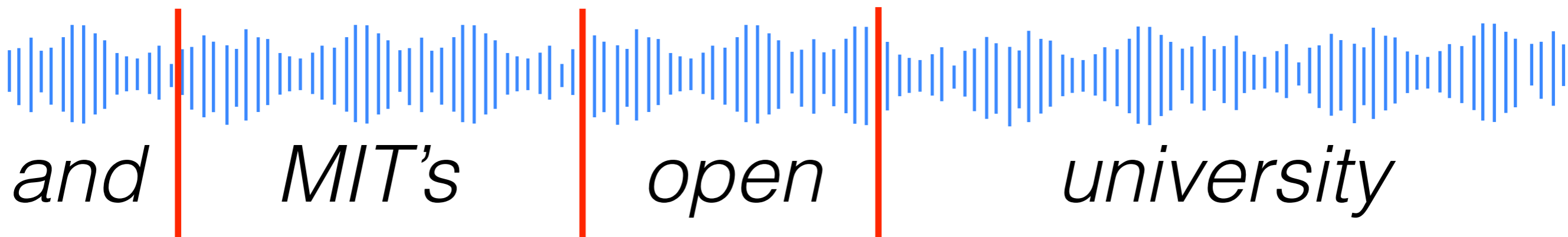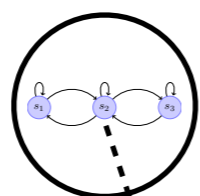## (Lee, O'Donnell, & Glass, 2015)



ɨ n d ɛ m a j tʰ ɪ z o w p ɨ ɲ j u w n ə v ɹ̩ z ɨ ɾ i

*and*    *MIT's*    *open*    *university*

# Unsupervised Lexicon Discovery from Speech
(Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

2. **Underlying and Surface Phones**

- Highly simplified phonetics/phonology.
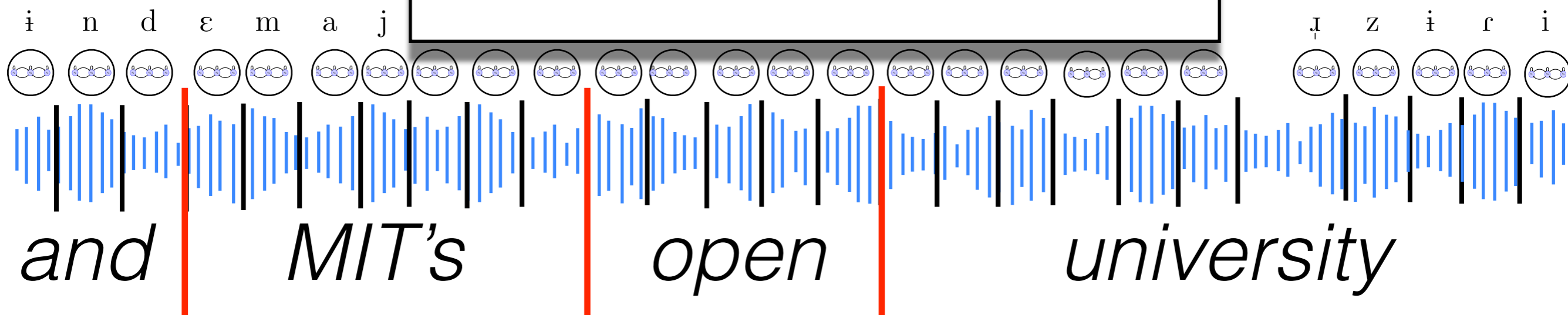  - Substitutions.
  - Splits.
  - Deletions.

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
### (Lee, O'Donnell, & Glass, 2015)



**3. Model of Morphological Structure**

- Utterances consist of words.
- Words consist of sub-words.
- Sub-words consist of phonemes.

*and*    *MIT's*    *open*    *university*

# Unsupervised Lexicon Discovery from Speech
### (Lee, O'Donnell, & Glass, 2015)

**3. Model of Morphological Structure**

$$\mathtt{U} \longrightarrow \mathtt{W}^*$$

$$\mathtt{W} \longrightarrow \mathtt{Sw}^*$$

$$\mathtt{Sw} \longrightarrow \mathtt{P}^*$$

$$\mathtt{P} \longrightarrow l$$

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
### (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)
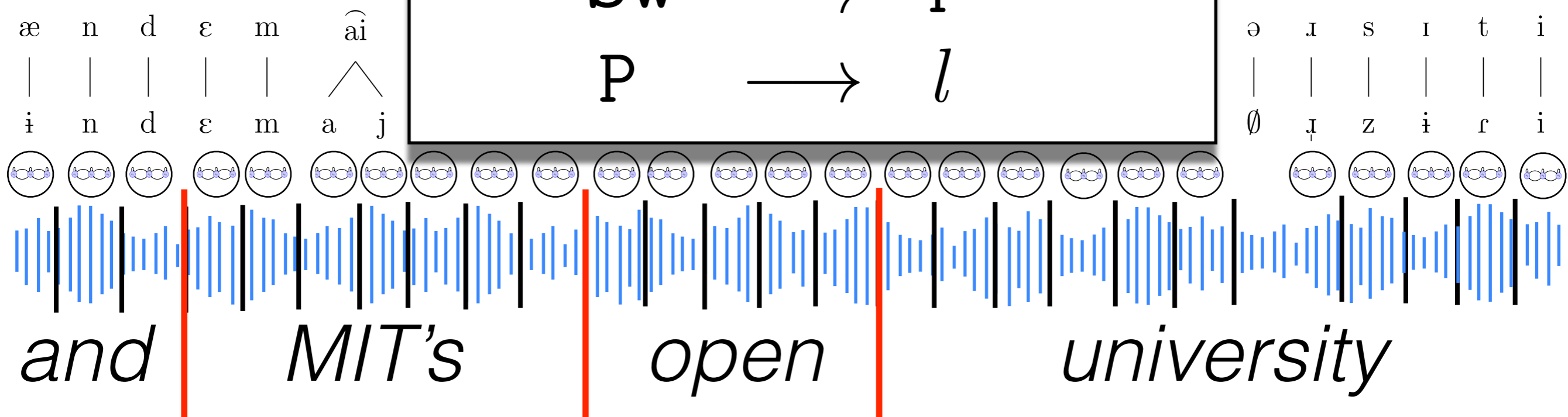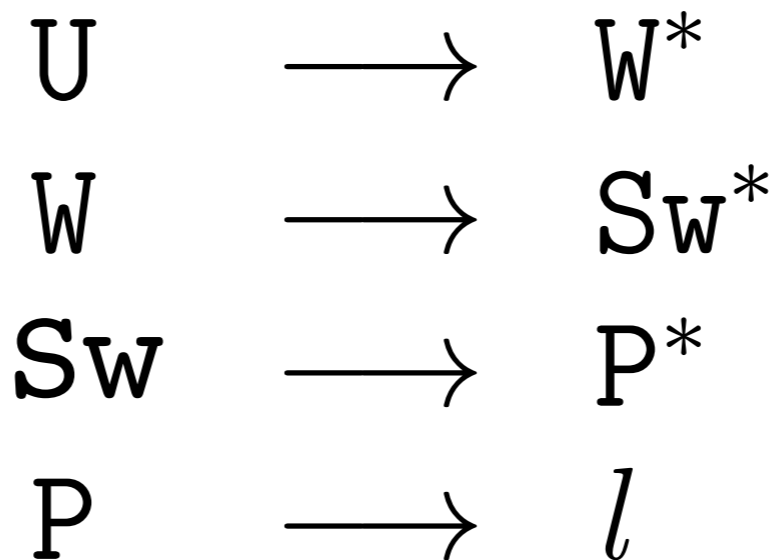


## 4. Model of Lexical Storage

- *Full-Storage* Model (Adaptor Grammars; Johnson, et al.. 2007)
  - First pass don't infer productivity.
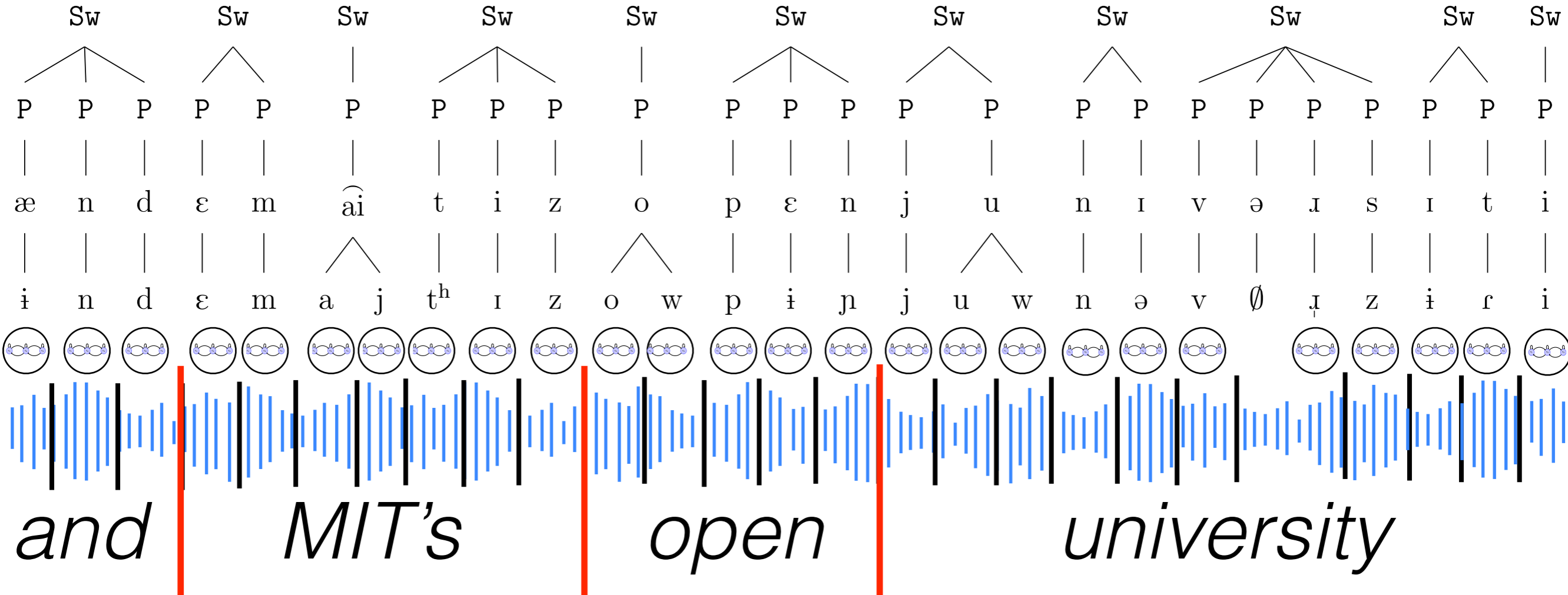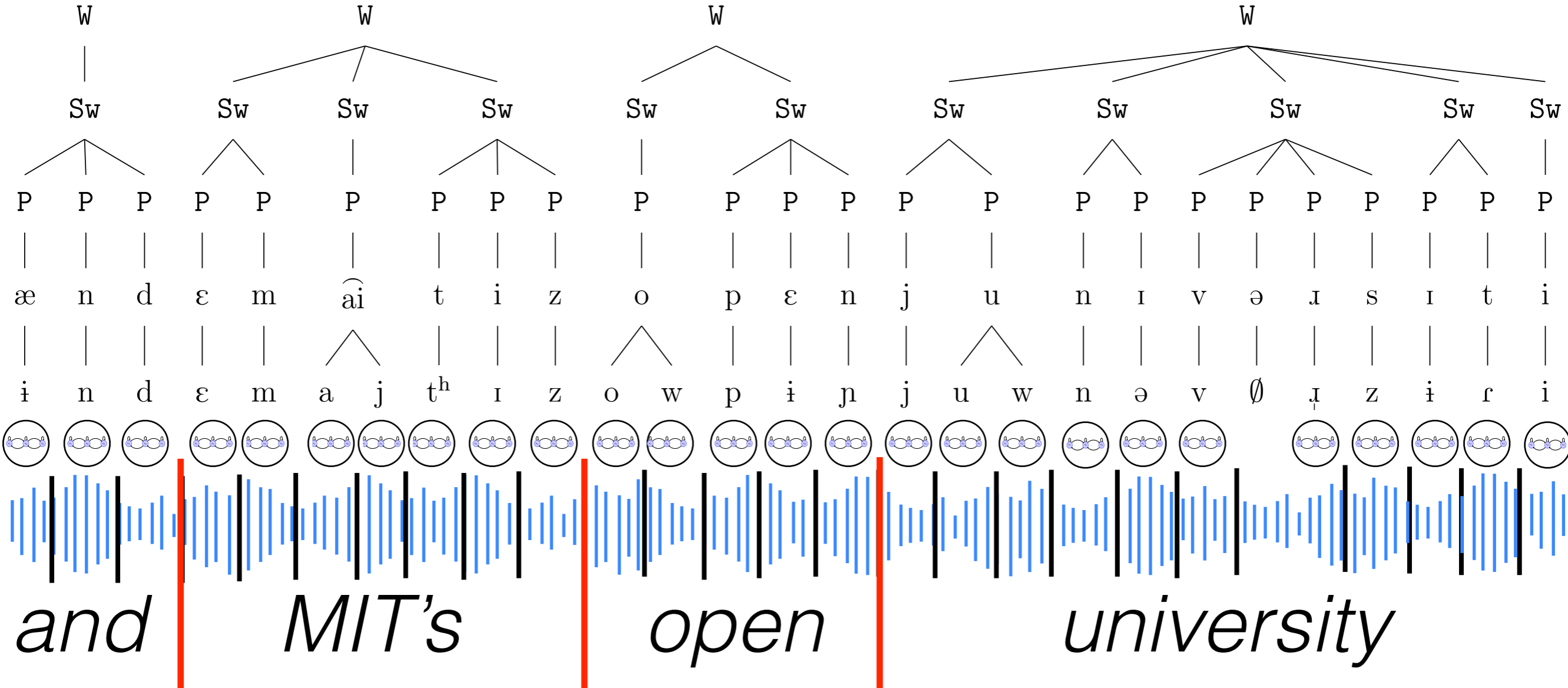- Store all sub-word and word units.

# Unsupervised Lexicon Discovery from Speech
(Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
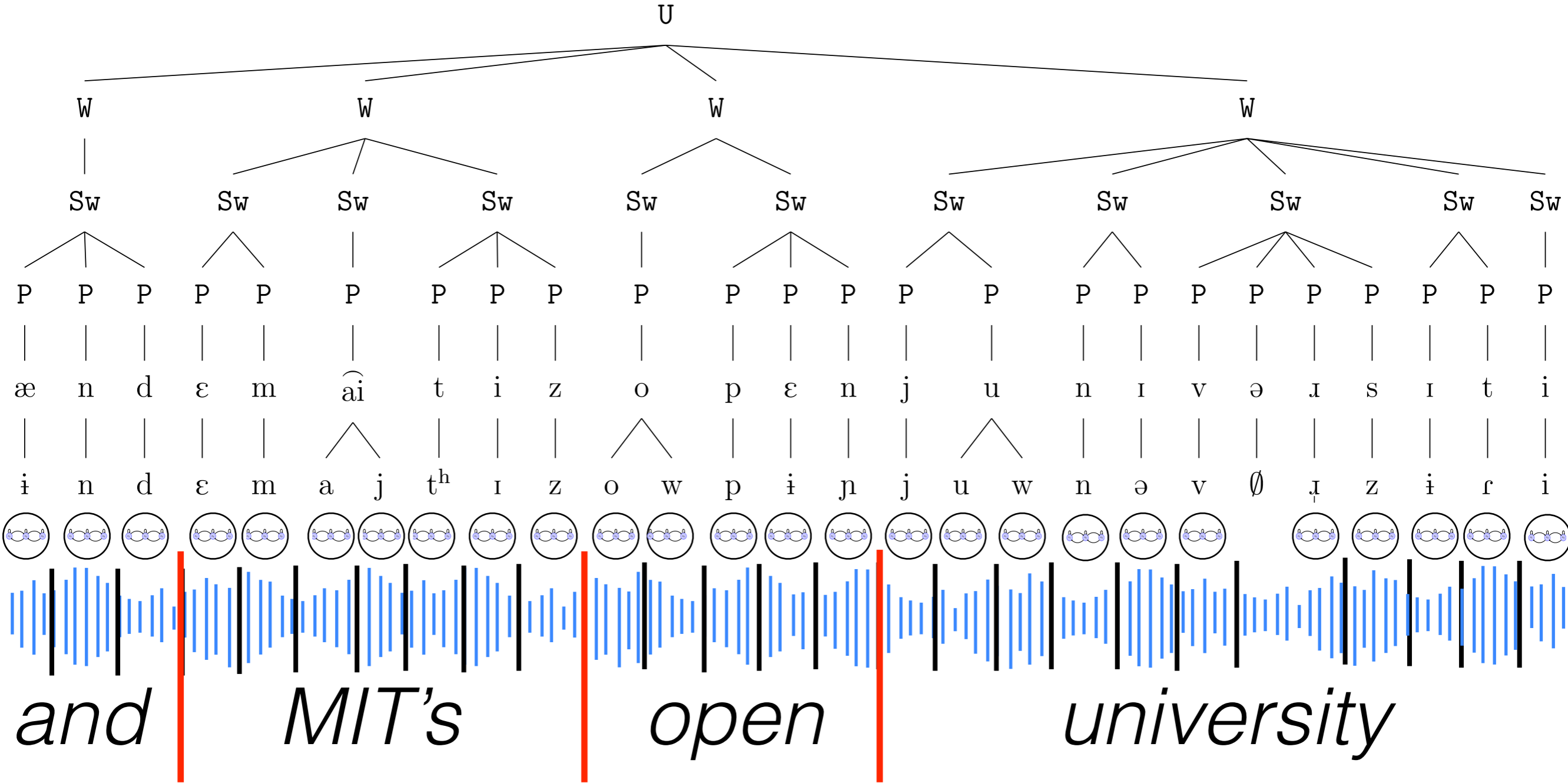## (Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech

(Lee, O'Donnell, & Glass, 2015)

# Unsupervised Lexicon Discovery from Speech
## (Lee, O'Donnell, & Glass, 2015)
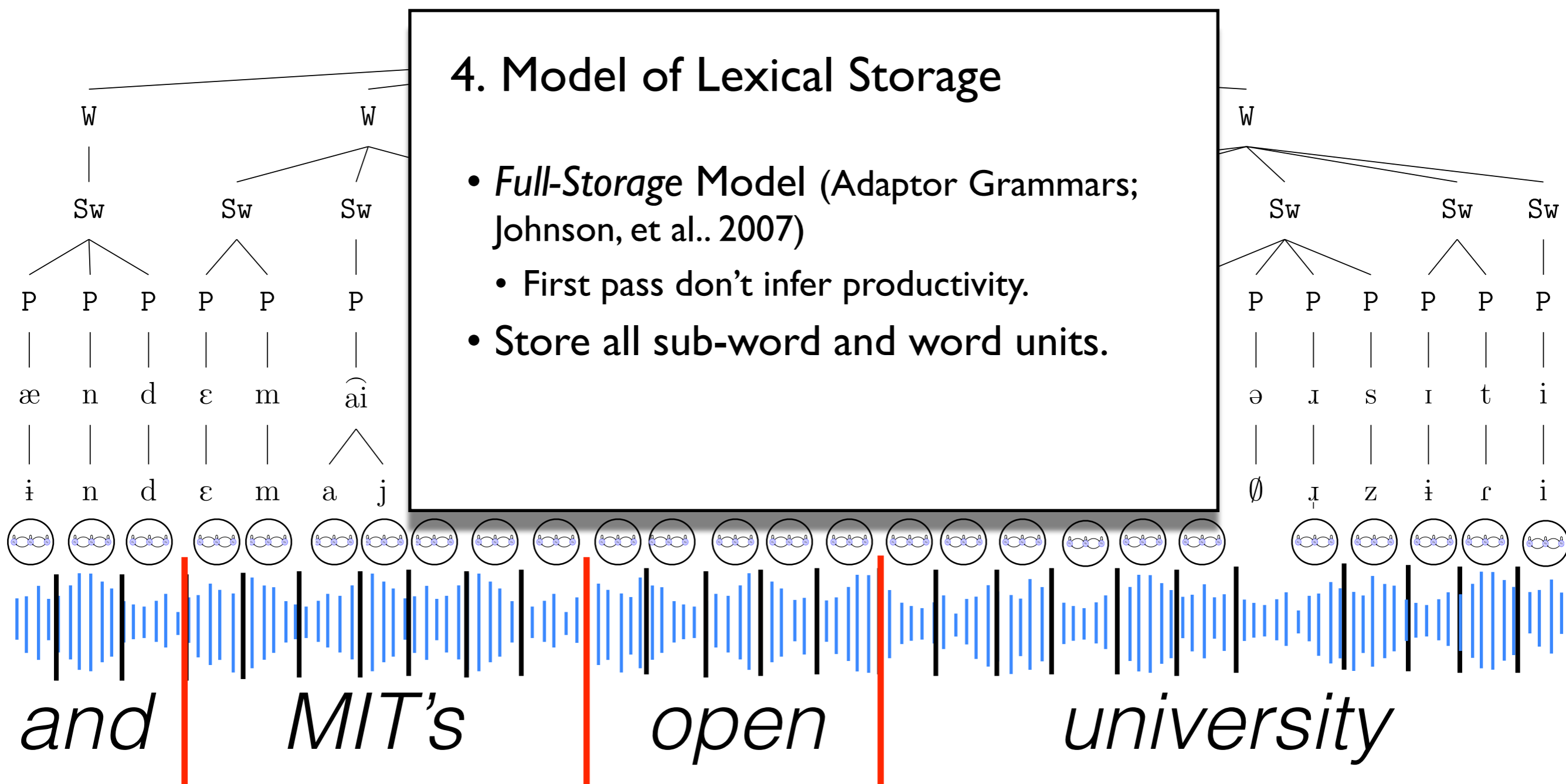
# Key Scientific Results

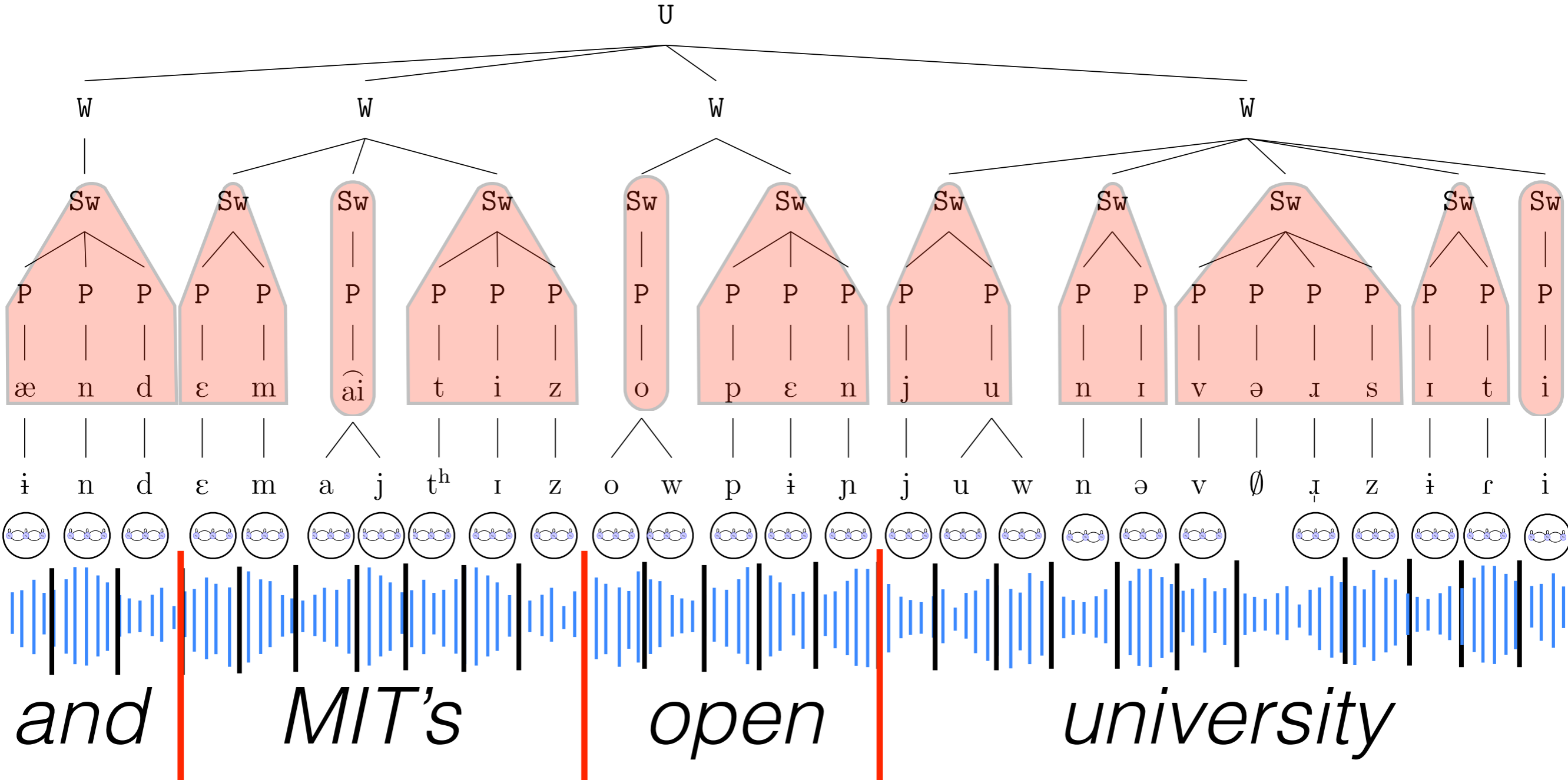- First unsupervised model of this scope.

- Outperforms state-of-art system for spoken term detection (Zhang, et al. 2013).

- Phone segmentation ~75%, word segmentation less great <20%.

- <u>Lexical units.</u>

- <u>Synergistic Interactions</u>.

# Key Scientific Results

- <u>Lexical units.</u>

- <u>Synergistic Interactions</u>.

# Lexical Unit Study

# Lexical Units

| Transcription | Discovered lexical units | |Word| |
|---|---|---|
| /iy l iy/ (really, willy, billion) | [35] [31 4]    -/ili/- | 68 |
| /ey sh ax n/ (innovation, imagination) | [6 7 30] [49] | 43 |
| /ax bcl ax l/ (able, cable, incredible) | [34 18] [38 91] | 18 |
| discovered | [26] [70 110 3] [9 99] [31] | 9 |
| individual | [49 146] [34 99] [154] [54 7] [35 48] | 7 |
| powerful | [50 57 145] [145] [81 39 38] | 5 |
| open university | [48 91] [4 67] [25 8 99 29] [44 22] [103 4] | 4 |
| the arab muslim world | [28 32] [41] [67] [25 35] [1 27] [13 173] [8 139] [38 91] | 2 |

# Lexical Units

| Transcription | Discovered lexical units | \|Word\| |
|---|---|---|
| /iy l iy/ (really, willy, billion) | [35] [31 4] | 68 |
| /ey sh ax n/ (innovation, imagination) | [6 7 30] [49] | 43 |
| /ax bcl ax l/ (able, cable, incredible) | [34 18] [38 91] | 18 |
| discovered | [26] [70 110 3] [9 99] [31] | 9 |
| individual | [49 146] [34 99] [154] [54 7] [35 48] | 7 |
| powerful | [50 57 145] [145] [81 39 38] | 5 |
| open university | [48 91] [4 67] [25 8 99 29] [44 22] [103 4] | 4 |
| the arab muslim world | [28 32] [41] [67] [25 35] [1 27] [13 173] [8 139] [38 91] | 2 |

*-ation*

# Predicted Generalizable Combinations

| Sequence | Category |
| --- | --- |
| *-ate -ion* | N |
| *-ic -al* | A |
| *-ate -ive* | A |
| *-al -ity* | N |
| *-al -ize* | V |
| *-ology -ist* | N |
| *-ment -al* | A |
| *-able -ity* | N |
| *-ist -ic* | A |
| *-ous -ity* | N |

# Lexical Units

| Transcription | Discovered lexical units | \|Word\| |
|---|---|---|
| /iy l iy/ (really, willy, billion) | [35] [31 4] | 68 |
| /ey sh ax n/ (innovation, imagination) | [6 7 30] [49] | 43 |
| /ax bcl ax l/ (able, cable, incredible) | [34 18] [38 91] *-able* | 18 |
| discovered | [26] [70 110 3] [9 99] [31] | 9 |
| individual | [49 146] [34 99] [154] [54 7] [35 48] | 7 |
| powerful | [50 57 145] [145] [81 39 38] | 5 |
| open university | [48 91] [4 67] [25 8 99 29] [44 22] [103 4] | 4 |
| the arab muslim world | [28 32] [41] [67] [25 35] [1 27] [13 173] [8 139] [38 91] | 2 |

# Lexical Units

| Transcription | Discovered lexical units | \|Word\| |
|---|---|---|
| /iy l iy/ (really, willy, billion) | [35] [31 4] | 68 |
| /ey sh ax n/ (innovation, imagination) | [6 7 30] [49] | 43 |
| /ax bcl ax l/ (able, cable, incredible) | [34 18] [38 91] | 18 |
| discovered | [26] [70 110 3] [9 99] [31] | 9 |
| individual | [49 146] [34 99] [154] [54 7] [35 48] | 7 |
| powerful | [50 57 145] [145] [81 30 38] | 5 |
| open university | [48 91] [4 ... | 4 |
| the arab muslim world | [28 32] [41] [67] [25 35] [1 27] [13 173] [8 139] [38 91] | 2 |

*The arab muslim world.*

# Key Scientific Results

- <u>Lexical units.</u>

- <u>Synergistic Interactions</u>.

# Synergistic Interactions

Integrating multiple components often makes learning better. of system mutually constrain one linguistic structures they learn are more accurate than if they were learned independently.

- <u>Lexical units and Phones:</u> Top-down influence of lexical unit learning improves phone unit learning.

- <u>Underlying-to-Surface Mapping and Words</u>: Modeling underlying-to-surface mapping improves ability to find correct words and sub-words.
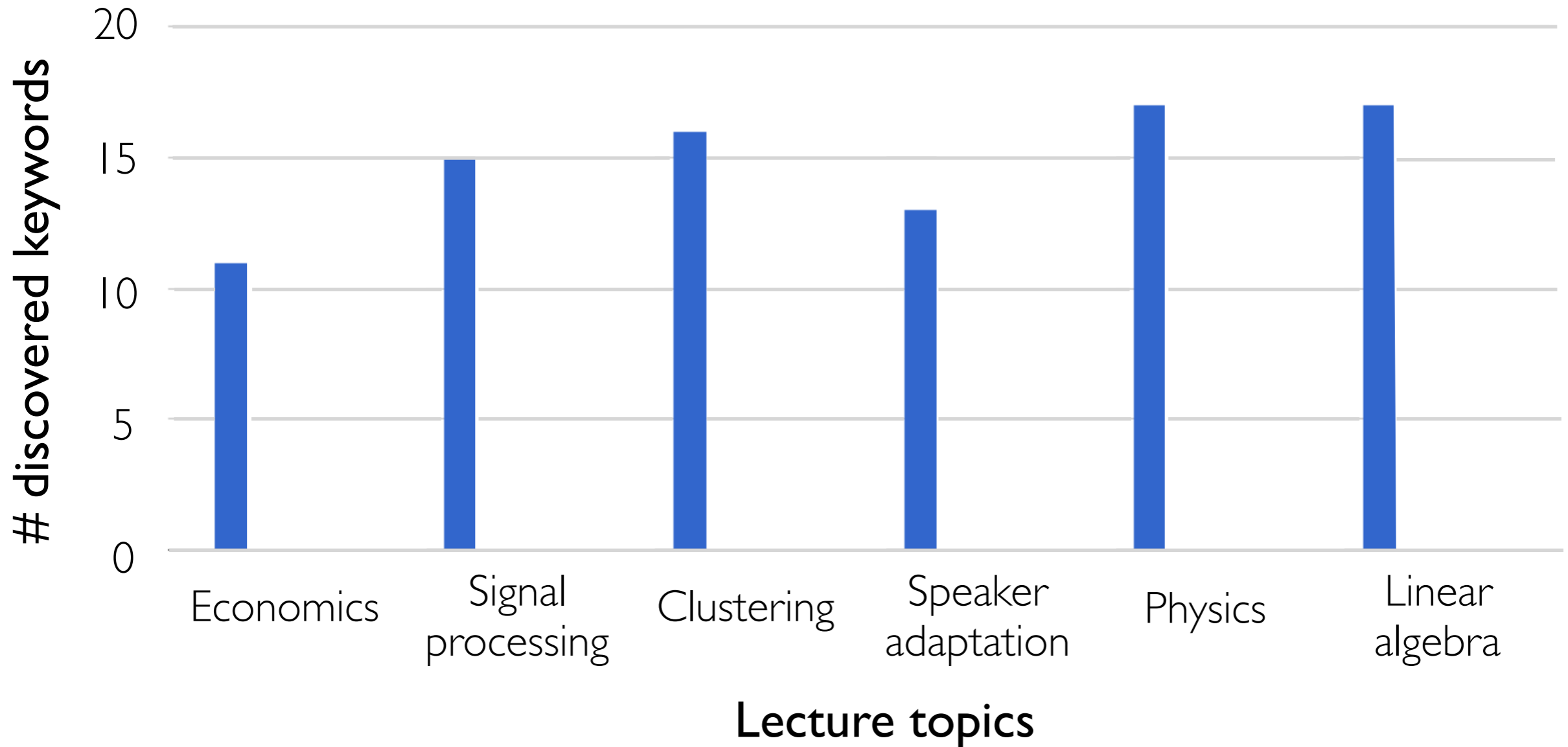
# Lesion Study

- **Lesioned** model components to study importance of difference components.

  - Examined top 20 most important words in each lecture (term frequency inverse-document frequency)
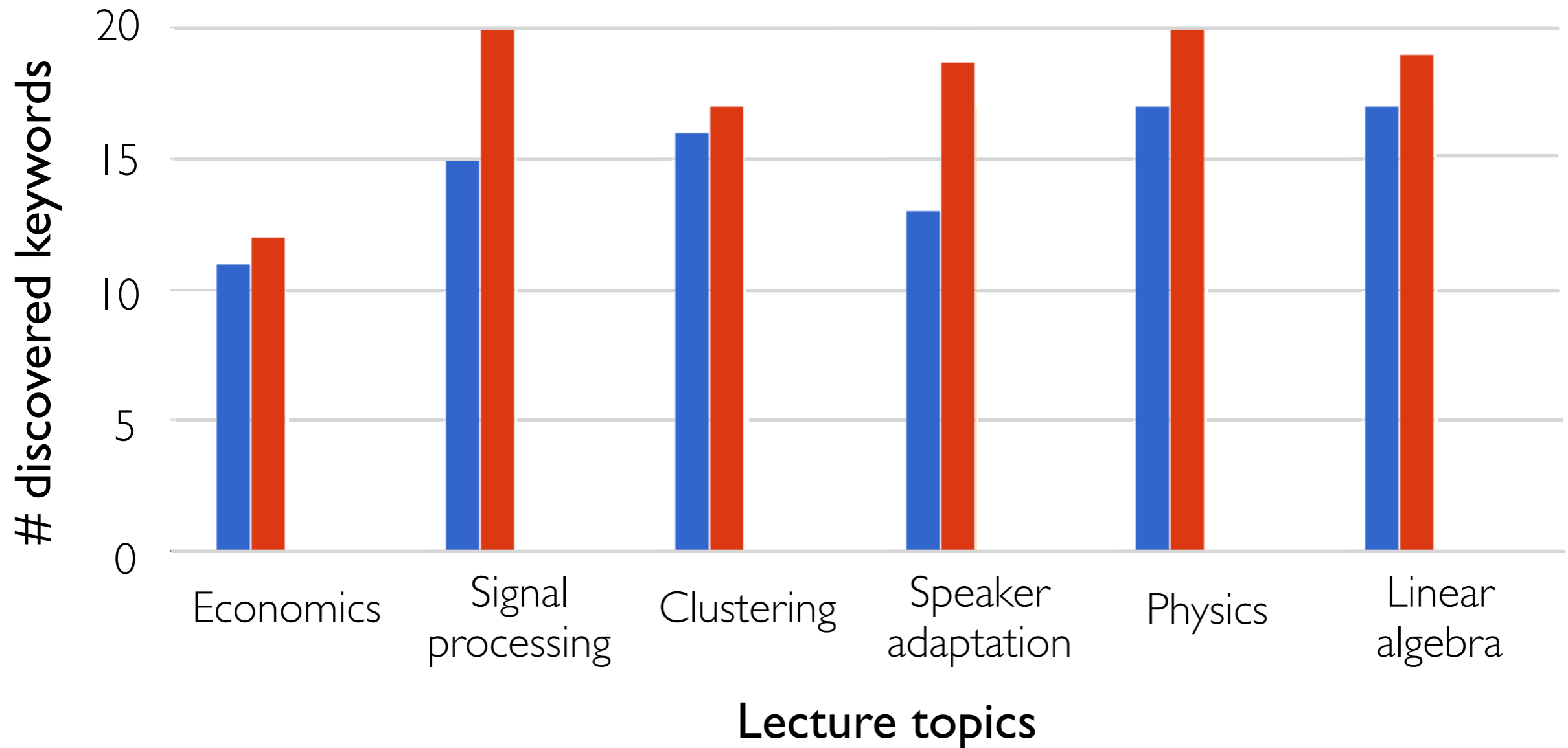
# Baseline

# Fixed Phones

# Joint Model

# No Noisy Channel

# No Noisy-Channel

# Overall Conclusion

- Can use tradeoff-based approach to learn lexical units.

- Makes fine-grained predictions for linguistics and psychology.

- Can scale up to more unsupervised settings.

# Thanks