

# ITERATIVE BAYESIAN WORD SEGMENTATION FOR UNSUPERVISED VOCABULARY DISCOVERY FROM PHONEME LATTICES

Jahn Heymann, Oliver Walter, Reinhold Häb-Umbach

University of Paderborn, Germany  
{heyman,walter,haeb}@nt.uni-paderborn.de  
http://nt.uni-paderborn.de

Bhiksha Raj

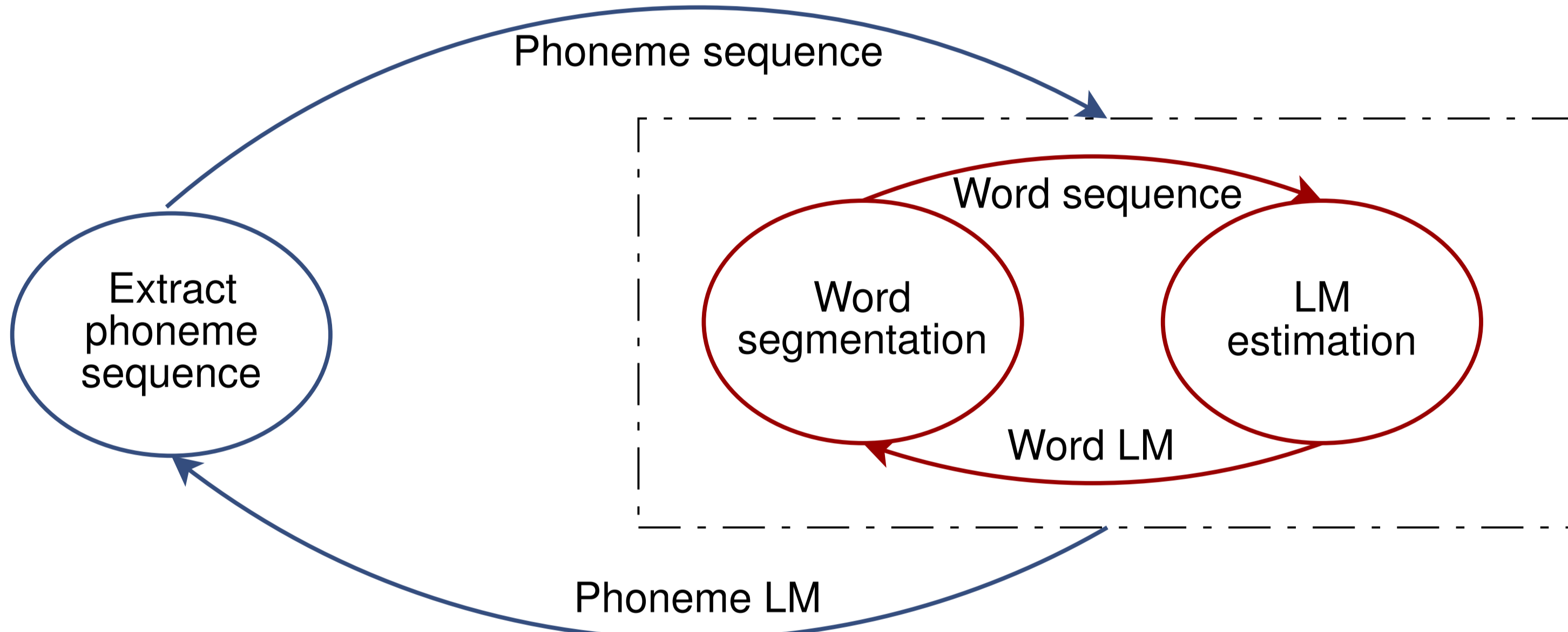
Carnegie Mellon University, USA  
bhiksha@cs.cmu.edu  
http://mlsp.cs.cmu.edu

## Introduction

- **Framework:** zero-resource speech recognition
- **Input:** Phoneme lattice generated by ASR engine
- **Goal:** Discover vocabulary, segment phoneme strings
- **Approach:**
  - ▶ Exploit consistency of character sequence within words
  - ▶ Simultaneous word segmentation and LM estimation

## Iterative 2-step Algorithm

- **Objective:** Maximize sentence likelihood
- **Iterate:** 1-best sequence extraction and word segmentation



- Learn two *n*-gram language models
  1. Word based *nested* PYLM for segmentation
  2. Phoneme based *hierarchical* PYLM for extraction
    - ▶ incorporates segmentation!
- Increase the order of the models after  $k_{sw}$  iterations
  - ▶ Start with a low model order for a initialization
  - ▶ Switch to a higher order for fine-tuning
- WFST-based implementation [Heymann13]

## Pitman-Yor Language Model [Teh06]

- Non-parametric i.e. unknown number of words
- Bayesian approach with power law prior (Zipf's law)
- Probability for word  $w$  in context  $\mathbf{u}$  recursively calculated as

$$\Pr(w|\mathbf{u}, S, \Theta) = \frac{c_{uw} \cdot \theta_{|\mathbf{u}|} + d_{|\mathbf{u}|} t_{uw}}{\theta_{|\mathbf{u}|} + c_{u..} + d_{|\mathbf{u}|} t_u} \Pr(w|\pi(\mathbf{u}), S, \Theta)$$

- Nesting: For new word use likelihood of word being character (phone) sequence  $c_1, \dots, c_k$  (fall back):

$$\Pr(w_{new}) \sim \prod_{i=1}^k \Pr(c_i | c_{i-1}, \dots, c_1, S, \Theta)$$

- Probability for characters (phones) calculated as above

## References

- [Neubig10] Learning a Language Model from Continuous Speech: G. Neubig, M. Mimura, S. Mori, T. Kawahara, InterSpeech 2010  
[Teh06] A hierarchical Bayesian language model based on Pitman-Yor processes: YW. Teh, ACL 2006  
[Heymann13] Unsupervised word segmentation from noisy input: J. Heymann, O. Walter, R. Häb-Umbach, B. Raj, ASRU 2013

## Experimental Setup

- **Lattice generation:** Monophone recognition using HTK
- **Dataset:** WSJCAM0 training set
  - ▶ 5628 sentences ( $\approx 10$  h speech), 10k vocabulary
  - ▶ Initial PER on best path: 33%
  - ▶ Apply different pruning factors to control lattice density
  - ▶ Minimal PER within lattice depends on pruning ( $\rightarrow L$ -PER)
- **Configurations:** (n-gram orders)

	word	phoneme I	phoneme II	algorithm
● (1)	1 → 2	2 → 8	4 → 8	proposed
▼ (2)	1	2	4	proposed
★ (3)	1	2	-	[Neubig10]
◆ (4)	2	8	-	[Neubig10]

## Experimental results

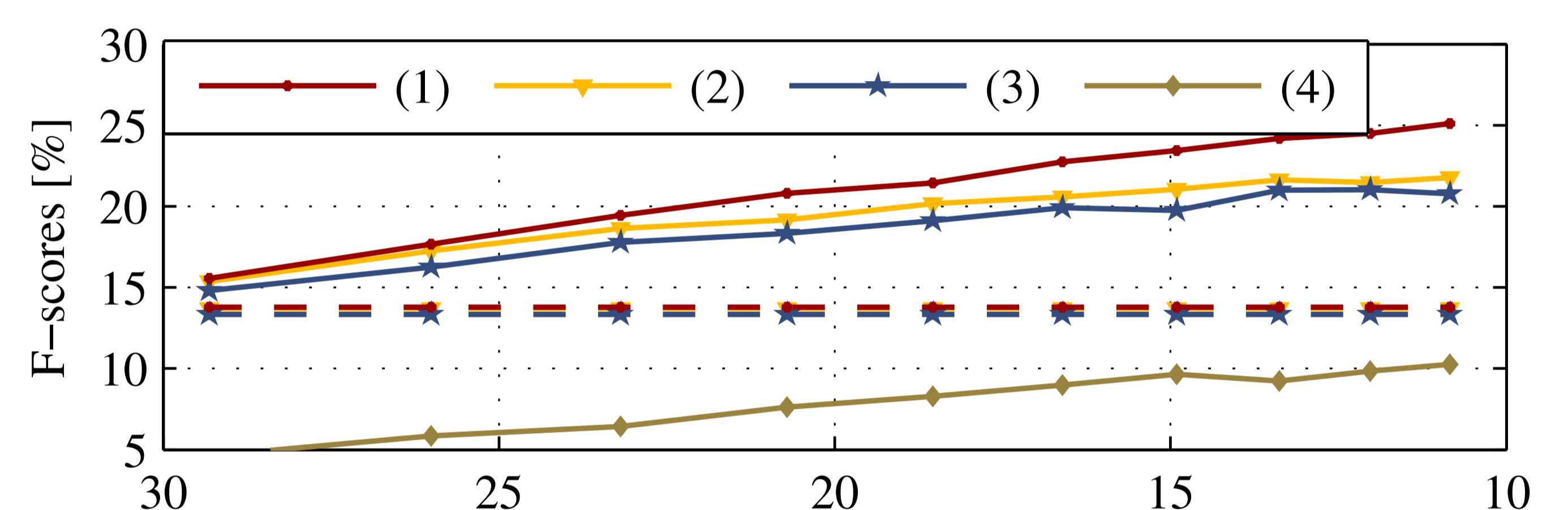


Figure 1: F-Score over L-PER for different setups. Dashed: result on best path

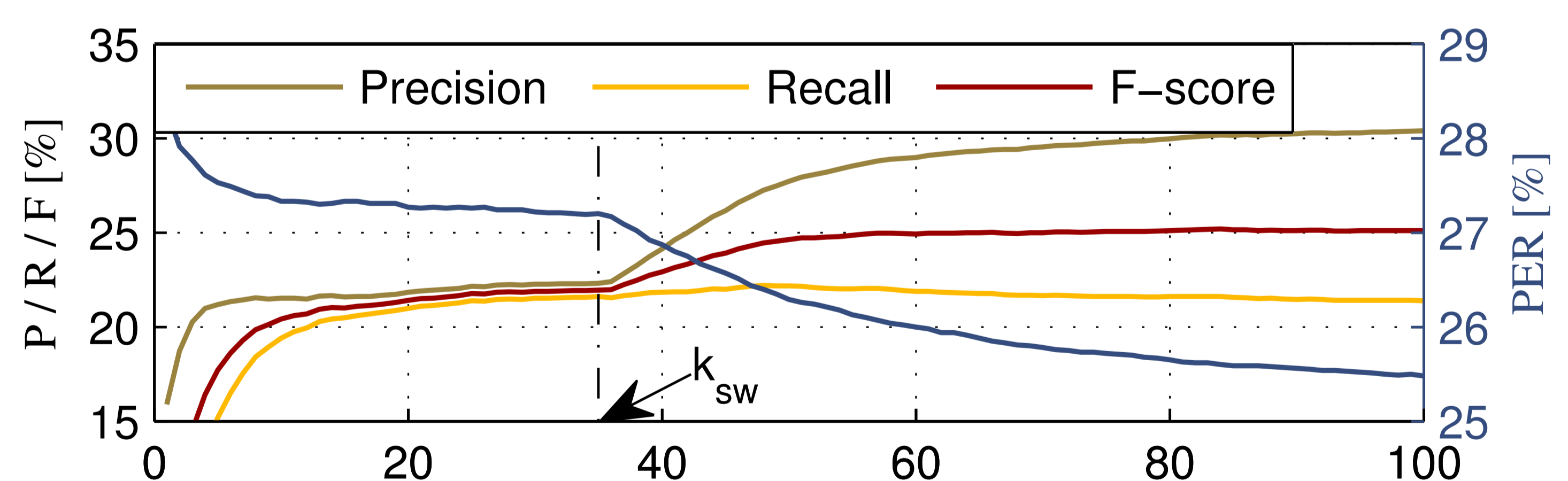


Figure 2: Measurements over iterations when switching the model order at  $k_{sw} = 35$ .



Figure 3: Top 100 discovered phonetic words (70 correct words)

## Conclusions

- Most words occurring more than ten times are found
- PER reduces from 33% to 25.5% by discovering vocabulary
- **Outlook:** Combine with unsupervised acoustic modeling

