

# Testing the robustness of online word segmentation: Effects of linguistic diversity and phonetic variation

Luc Boruta<sup>1,2</sup>, Sharon Peperkamp<sup>2</sup>, Benoît Crabbé<sup>1</sup> & Emmanuel Dupoux<sup>2</sup>

luc.boruta@inria.fr

<sup>1</sup>ALPAGE, Univ. Paris 7 & INRIA

<sup>2</sup>LSCP-DEC, EHESS, ENS & CNRS

CMCL — June 23, 2011

# Yet another study on word segmentation...

What this work is not about

- New models of word segmentation.

What this work is about

- The acquisition of **word segmentation**;
- The acquisition of **phonological knowledge**;
- Interactions between the two.

# Yet another study on word segmentation...

## What this work is not about

- New models of word segmentation.

## What this work is about

- The acquisition of **word segmentation**;
- The acquisition of **phonological knowledge**;
- Interactions between the two.

## Yet another study on word segmentation...

### What this work is not about

- New models of word segmentation.

### What this work is about

- The acquisition of **word segmentation**;
- The acquisition of **phonological knowledge**;
- Interactions between the two.

# Word segmentation vs. allophonic rules

## French devoicing allophonic rule

$$/r/ \rightarrow \begin{cases} [x] & \text{before a voiceless consonant} \\ [ʁ] & \text{otherwise} \end{cases}$$

## Consequence

$$/kanar/ \rightarrow \begin{cases} [kana_x\_flotã], \textit{canard flottant} \\ [kana_ʁ\_ʒon], \textit{canard jaune} \end{cases}$$

# Word segmentation vs. allophonic rules

## French devoicing allophonic rule

$$/r/ \rightarrow \begin{cases} [x] & \text{before a voiceless consonant} \\ [ʁ] & \text{otherwise} \end{cases}$$

## Consequence

$$/kanar/ \rightarrow \begin{cases} [kana_x\text{̣}flotã], \textit{canard flottant} \\ [kana_ʁ\text{̣}ʒon], \textit{canard jaune} \end{cases}$$

# Word segmentation

## The task

- Input: /əwʊdʃlɪkwʊdʃlɪkwʊd/
- Output: /ə\_wʊdʃlɪk\_wʊd\_ʃlɪk\_wʊd/

## Phonemic transcripts = idealized input

- Models are typically evaluated using **phonemic** transcripts;
- Assumption: kids know how to undo **allophony**/coarticulation.

# Word segmentation

## The task

- Input: /əwʊdʃlɪkwʊdʃlɪkwʊd/
- Output: /ə\_wʊdʃlɪk\_wʊd\_ʃlɪk\_wʊd/

## Phonemic transcripts = idealized input

- Models are typically evaluated using **phonemic** transcripts;
- Assumption: kids know how to undo **allophony**/coarticulation.



## Related work

### Rytting, Brew & Fosler-Lussier (2010)

- Input unit: probability vector over a finite set of symbols;
- Symbols: limited to the phonemic inventory.

### Daland & Pierrehumbert (2010)

- Input: phonemic transcripts, conversational reduction processes;
- Reduction processes: implemented by hand;
- Transcripts: adult-directed speech.

## Related work

### Rytting, Brew & Fosler-Lussier (2010)

- Input unit: probability vector over a finite set of symbols;
- Symbols: limited to the phonemic inventory.

### Daland & Pierrehumbert (2010)

- Input: phonemic transcripts, conversational reduction processes;
- Reduction processes: implemented by hand;
- Transcripts: adult-directed speech.

# Which segmentation models?

## Desirable properties

[Brent, 1999; Gambell & Yang, 2004]

- Start without any knowledge specific to a particular language;
- Learn in an unsupervised manner and operate incrementally.

## Which segmentation models?

- MDBP-1: Brent, 1999;
- NGS-u: Venkataraman, 2001;
- Two random baselines.

# Which segmentation models?

## Desirable properties

[Brent, 1999; Gambell & Yang, 2004]

- Start without any knowledge specific to a particular language;
- Learn in an unsupervised manner and operate incrementally.

## Which segmentation models?

- MDBP-1: Brent, 1999;
- NGS-u: Venkataraman, 2001;
- Two random baselines.

# Evaluation

## Now-standard evaluation protocol

[Brent, 1999; Goldwater et al., 2009]

- Gold standard: orthographic segmentation;
- Precision, recall and **F-score** on the word segmentation;
- Precision, recall and **F-score** on the induced lexicon.

	Lexicon	Segmentation
ə_ʷsɔtʃlɔk_ʷsɔd_tʃlɔk_ʷsɔd	✓	✓
ə_ʷsɔd_tʃlɔk_ʷsɔd_tʃlɔk_ʷsɔd	✓	✗

## Experimental setup

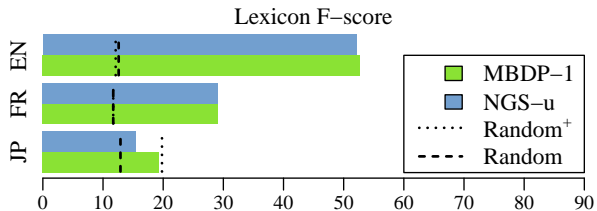
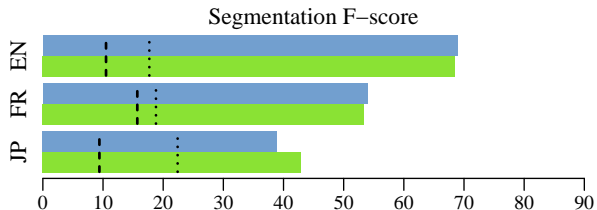
### CHILDES corpora of child-directed speech

[MacWhinney, 2000]

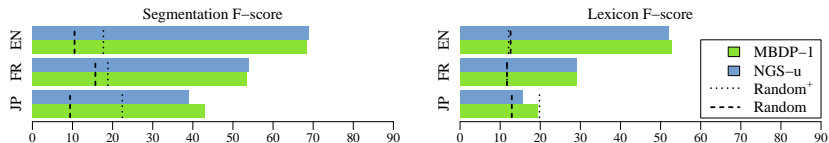
- Derived from transcribed adult-child verbal interactions;
- Phonemic transcriptions, orthographic segmentation.

	English	French	Japanese
Utterance tokens	10k	10k	10k
Word tokens	33k	51k	27k
Phoneme tokens	96k	121k	103k
Phoneme types	50	35	49

# Cross-linguistic evaluation on phonemic corpora



# Cross-linguistic evaluation on phonemic corpora



- Blame it on the data?
- Rich morphology (e.g. French clitics)? Hapax rate?
- Relative importance of different cues?



# Effects of phonetic variation

## Phonemic transcripts = idealized input

- Models are typically evaluated using **phonemic** transcripts;
- Assumption: kids know how to undo **allophony**/coarticulation.

## Corpora and allophonic rules

- No phonetic transcripts of child-directed speech are available;
- How many allophones do infants have to learn?
- Where is the limit between allophony and mere coarticulation?

# Effects of phonetic variation

## Phonemic transcripts = idealized input

- Models are typically evaluated using **phonemic** transcripts;
- Assumption: kids know how to undo **allophony**/coarticulation.

## Corpora and allophonic rules

- No phonetic transcripts of child-directed speech are available;
- How many allophones do infants have to learn?
- Where is the limit between allophony and mere coarticulation?

# Experimental setup

## Emulating rich phonetic transcriptions

[Boruta, 2011a]

- Apply **artificial allophonic rules** to phonemic corpora;
- Benchmark models using different **allophonic complexities**;
- Control the size of the allophonic grammar.

## Simplifying assumptions

[LeCalvez, 2007; Boruta, 2011a]

- We only model monolateral rules:  $p \rightarrow a / \_ c$
- No two rules introduce the same phone:  $[r]^{/t/}$  and  $[r]^{/d/}$

# Experimental setup

## Emulating rich phonetic transcriptions

[Boruta, 2011a]

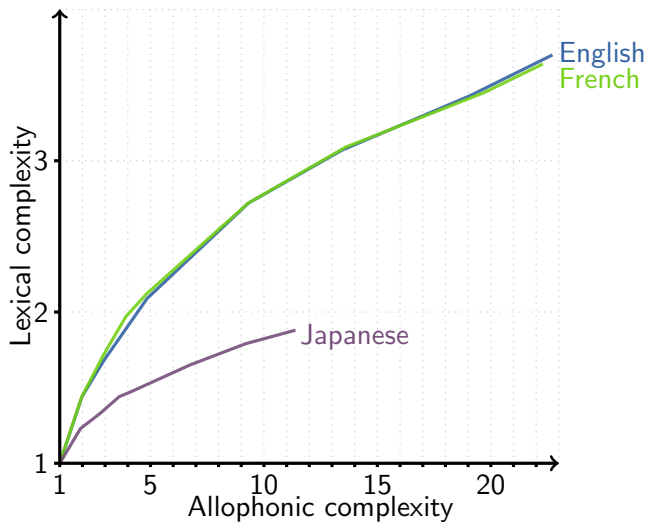
- Apply **artificial allophonic rules** to phonemic corpora;
- Benchmark models using different **allophonic complexities**;
- Control the size of the allophonic grammar.

## Simplifying assumptions

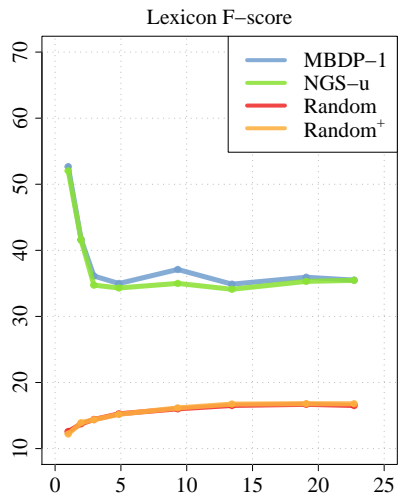
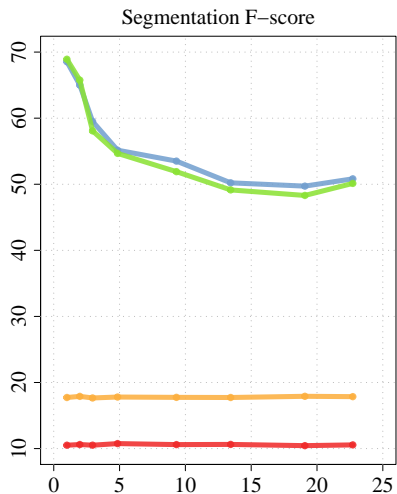
[LeCalvez, 2007; Boruta, 2011a]

- We only model monolateral rules:  $p \rightarrow a / \_ c$
- No two rules introduce the same phone:  $[r]^{/t/}$  and  $[r]^{/d/}$

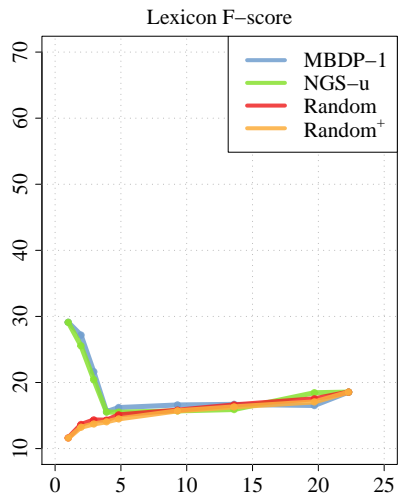
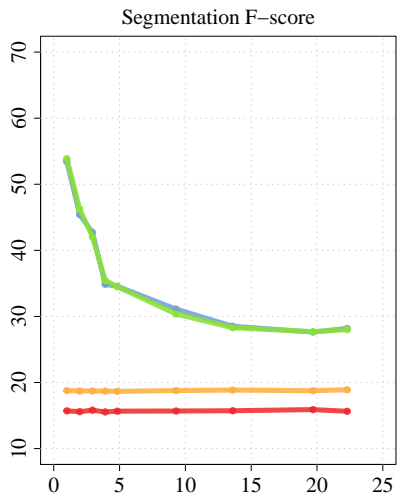
## Lexical complexity $\propto$ allophonic complexity



# Results: English

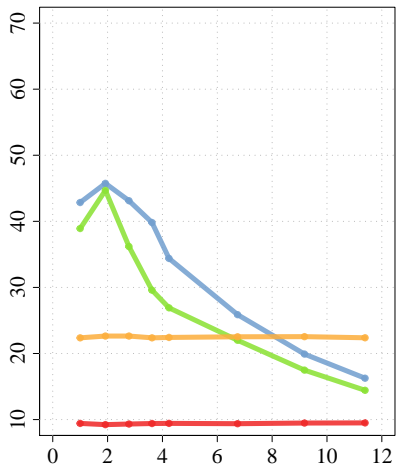


# Results: French

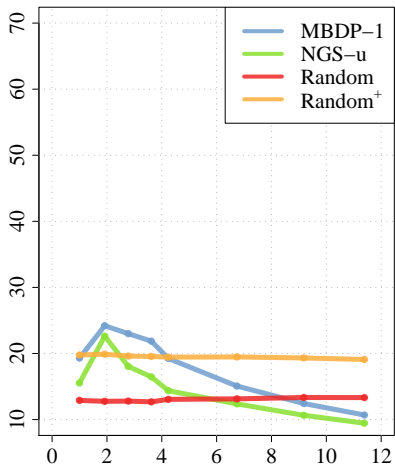


# Results: Japanese

## Segmentation F-score

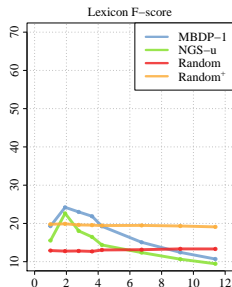
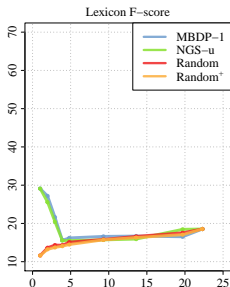
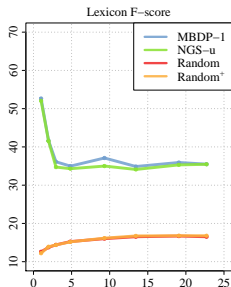


## Lexicon F-score





# Effects of phonetic variation



## Unsurprising results

- No mechanism for 'explaining away' allophonic variation;
- Any word form found by the models will be added to the lexicon.

# Conclusion

## Take-home message

- Cross-linguistic evaluation is not dispensable;
- Phonetic inputs impact word seg. models' performance.
- Phonological knowledge < word segmentation?

## Where to go from here?

- Incorporate some mechanism to handle noise and/or variation;
- Use the imperfect lexical knowledge to help learning a phonology.

↪ *Combining indicators of allophony*, ACL'11 Student Session.

# Conclusion

## Take-home message

- Cross-linguistic evaluation is not dispensable;
- Phonetic inputs impact word seg. models' performance.
- Phonological knowledge < word segmentation?

## Where to go from here?

- Incorporate some mechanism to handle noise and/or variation;
- Use the imperfect lexical knowledge to help learning a phonology.

↪ *Combining indicators of allophony*, ACL'11 Student Session.

/θæŋk⌵ju: /