#### The Grammar(s) of Code-Mixing: A corpus-based approach to multilingual first language acquisition



Antje Endesfelder Quick University of Leipzig





**Stefan Hartmann** University of Düsseldorf

Nikolas Koch LMU Munich

## Code-Mixing

und dann magic air 'And then magic air'

Look at the Ampel, it's kaput 'Look at the traffic light, it's broken'

We have to beeil 'We have to hurry up' Code-Mixing = one of the more salient phenomena found in bilingualism → use of two languages in one utterance

Recent years have seen increased interest in codemixing from a usage-based perspective

#### How do children acquire language(s)

#### Complexity of languages is remarkable - Different research traditions





 ✓ children are equipped with pre-existing experience (e.g. Chomsky 1965)  ✓ children acquire languages by actively constructing complexity (e.g. Tomasello 2003) → piecemeal acquisition



#### Building up language (s)



Blurring the line between lexicon and grammar



• Hypothesis: language acquisition is strongly itembased – early child language is highly formulaic



#### How to account for patterns ?



### The Traceback method



## The Traceback method



## The Traceback method



#### How to account for patterns ?



### Chunk-Based Learner (CBL)



## Chunk-Based Learner (CBL)



detecting chunks incrementally → recognizing multiword chunks by using backward transitional probabilities (BTP)



### CBL – calculating BTP





### CBL – calculating BTP



## CBL – identifying chunks



## CBL – identifying chunks

p(the | stray) > AVG.TP p(stray | dog) > AVG.TP





# Data - Bilingual Corpus



 ✓ entire dataset → monolingual utterances and input as input for the CBL algorithm
✓ focus on the "comprehension" side here, i.e., the chunks that the model identifies

#### Language Proportions



script for the CBL algorithms available @ https://github.com/StewartMcCauley/CBL/

#### **Results - Traceback**



#### **Results - Traceback**





#### **Results CMed - CBL**





✓ lots of CMed utterances contain a chunk or a partial chunk
✓ differences between the chunks identified in the child's code-mixed and his monolingual utterances → model does not make a difference between them (G-E-CM)

#### **Results CMed - CBL**



- ✓ position of chunk boundaries identified by the CBL algorithm often coincides with the position of code-mixes
  - ✓ because an English word is of course relatively unlikely to be preceded by a German one, and vice versa
- *i. ein kleinen* | *shark* (a little shark)
- *ii. nein* || a || *nein* || a *ice hockey player* (no a no a ice hockey player)
- iii. zeig | | ice cream (show ice cream)

- no method has been proven to be exhaustive and that depends on the ways how they implement the notion of patterns
- CBL results complement the TB results in a useful way: both models detect formulaic language use but make different predictions that we see in the data
  - TB  $\rightarrow$  frame-and-slot patterns
  - CBL  $\rightarrow$  chunks
- frequent code-mixed sequences like *nein this*, which were identified as fixed chunks by the TB method, are not identified as chunks by the CBL algorithm

 code-mixed utterances can be accounted for by lexically fixed patterns and emerging frame-and-slot patterns → generating/recycling (creative) utterances from 'bits and pieces' of already acquired constructions

 important to view code-mixing as a dynamic process instead of trying to find a 'one-fit-all' grammar to account for mixes across different populations of children and different languages



nuqneH!