

gramophone – A hybrid approach to grapheme-phoneme conversion

Kay-Michael Würzner, Bryan Jurish
{wuerzner,jurish}@bbaw.de

FSMNLP
Universität Düsseldorf
24th June 2015



Overview

The Task

- Finding the pronunciation of a word given its spelling

The Challenge: Ambiguity

- a phoneme may be realized by different characters
- a character may be represented by different phonemes

Our Approach: A combination of

- a hand-crafted rule set controlling segmentation and alignment,
- a conditional random field model for generating transcription candidates, and
- an N -gram language model for selecting the “best” grapheme-phoneme mapping



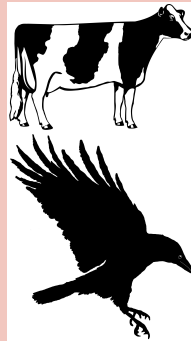
Outline

1. Grapheme-phoneme conversion and its applications
2. Existing approaches
3. The gramophone approach
 - (a) Alignment/Encoding
 - (b) Transcription
 - (c) Rating
4. Comparative evaluation and error analysis
5. Discussion & Outlook



Grapheme-phoneme conversion: Problem description

- Symbolic representation of the pronunciation of words
- Orthography is ambiguous w.r.t. pronunciation, phonetic alphabets allow for an unambiguous representation



cow /kaʊ/

crow /kɹoʊ/

- Complex alignment: Single characters may be represented by multiple phonemes (and *vice versa*)

ph oe n i x
f iː n ɪ ks



Grapheme-phoneme conversion: Applications

Text-to-speech systems

(Black & Taylor 1997)

- Improvement of speech signal synthesis by disambiguation of the input text

Spelling correction / “canonicalization”

(Jurish 2010)

- Phonetic transcriptions as a normal form for identifying spelling variants

Speech recognition

(Galescu and Allen 2002)

- Inverse application of g2p models

Pronunciation dictionaries

(TC-Star project; DWDS)

- Generation of transcriptions or transcriptions candidates especially in compounding languages



Previous work: Rule-based approaches

- Inspired by *The Sound Pattern of English* (Chomsky & Halle 1968)
- Equivalent to regular grammars and rewriting systems (Johnson 1972)
- Successful model for g2p converters in many languages
- Used in various text-to-speech systems, e.g.
 - MITalk (Allen et al. 1987)
 - TETOS (Wothke 1993)
 - festival (Taylor et al. 1998)
- **Drawbacks:**
 - Expertise and effort required in their production and maintenance
 - Treatment of exceptional pronunciation e.g. in loan words (or even worse *compounds* of foreign and native words)

Versaillesdiktat

/vɛʁzaɪdɪktat/

engl. 'Versailles diktat'



Previous work: Statistical approaches

- Automatic inference of regularities in the correspondence of spellings and pronunciations from data (i.e. word+transcription pairs)
- Many large data sets exist
 - NETTalk
 - CELEX
 - wiktionary
- Many more existing approaches
 - Neural networks
 - Joint-sequence N -gram models
 - Conditional random fields
- **Drawback:**
 - No direct control of results, linguistically implausible transcriptions may be inferred

(cf. Reichel et al. 2008)

(Sejnowski & Rosenberg 1987)

(Bisani & Ney 2008)

(Jiampojarn & Kondrak 2009)

Getue \mapsto */gəʈə/
engl. 'fuss'



Alignment

Starting point

- Association of transcriptions with entire words
 \rightsquigarrow *Alignment on the grapheme-substring level necessary*
- $n : m$ relation between grapheme-phoneme string pairs

$$n, m \in \mathbb{N} \setminus \{0\}$$



Alignment

Starting point

- Association of transcriptions with entire words
 \rightsquigarrow *Alignment on the grapheme-substring level necessary*
- $n : m$ relation between grapheme-phoneme string pairs

$$n, m \in \mathbb{N} \setminus \{0\}$$

ph	oe	n	i	x
↕	↕	↕	↕	↕
f	i:	n	i	ks



Alignment

Starting point

- Association of transcriptions with entire words
 \rightsquigarrow *Alignment on the grapheme-substring level necessary*
- $n : m$ relation between grapheme-phoneme string pairs

$$n, m \in \mathbb{N} \setminus \{0\}$$

Approaches

- Numerous existing alignment methods
- Simplify the $n : m$ relation to a more tractable case

(cf. Reichel 2012)

$$n, m \in \{0, 1\}$$

ph	oe	n	i	x
↕	↕	↕	↕	↕
f	i:	n	i	ks



Alignment

Starting point

- Association of transcriptions with entire words
 \rightsquigarrow *Alignment on the grapheme-substring level necessary*
- $n : m$ relation between grapheme-phoneme string pairs

$$n, m \in \mathbb{N} \setminus \{0\}$$

Approaches

- Numerous existing alignment methods
- Simplify the $n : m$ relation to a more tractable case

(cf. Reichel 2012)

$$n, m \in \{0, 1\}$$

p	h	o	e	n	i	x	ε
↕	↕	↕	↕	↕	↕	↕	↕
f	ε	i:	ε	n	i	k	s



Alignment

Starting point

- Association of transcriptions with entire words
 \rightsquigarrow *Alignment on the grapheme-substring level necessary*
- $n : m$ relation between grapheme-phoneme string pairs

$$n, m \in \mathbb{N} \setminus \{0\}$$

Approaches

- Numerous existing alignment methods
- Simplify the $n : m$ relation to a more tractable case
- Application of some Levenshtein-like mechanism

(cf. Reichel 2012)

$$n, m \in \{0, 1\}$$

(Levenshtein, 1966)

p	h	o	e	n	i	x	ε
↕	↕	↕	↕	↕	↕	↕	↕
f	ε	i:	ε	n	i	k	s



Alignment

Alternatives?

- *Deletion* doubtful in the context of grapheme-phoneme correspondence
- Inference of many-to-many alignments error-prone (*Jiampojarn et al. 2007*)
- Linguistically motivated alignment desirable

Constraint-based alignment

- **Manual** definition of possible mappings between grapheme sequences and phonemic realizations

$$M \subset (\Sigma_G^+ \times \Sigma_P^+)$$

- Compiled as FST

$$E = \langle Q, \Sigma_G \cup \{| \}, \Sigma_P \cup \{ - \}, q_0, q_0, \delta \rangle$$

- Add a path $(q_0, q_0, g \cdot |, p \cdot -)$ for each mapping $(g, p) \in M$
 - '|' and '-' are reserved delimiter symbols
- Generate **all admissible** segmentations of a word and its transcription
 - FST I_G with a path $(q_0, q_0, g, g \cdot |)$ for every g in the domain of M
 - FST I_P with a path $(q_0, q_0, p, p \cdot -)$ for every p in the codomain of M



Alignment

- Construct letter FSTs W and T for a word w and its transcription t
- Alignment of w and t is generated by a series of compositions which filters out all non-matching pairings

$$A_{W,T} = \pi_2(W \circ I_G) \circ E \circ \pi_2(T \circ I_P)$$

Example

$$M = \{u:/u/, u:/u\:/, u:/ju\:/, uu:/u\:/\}$$



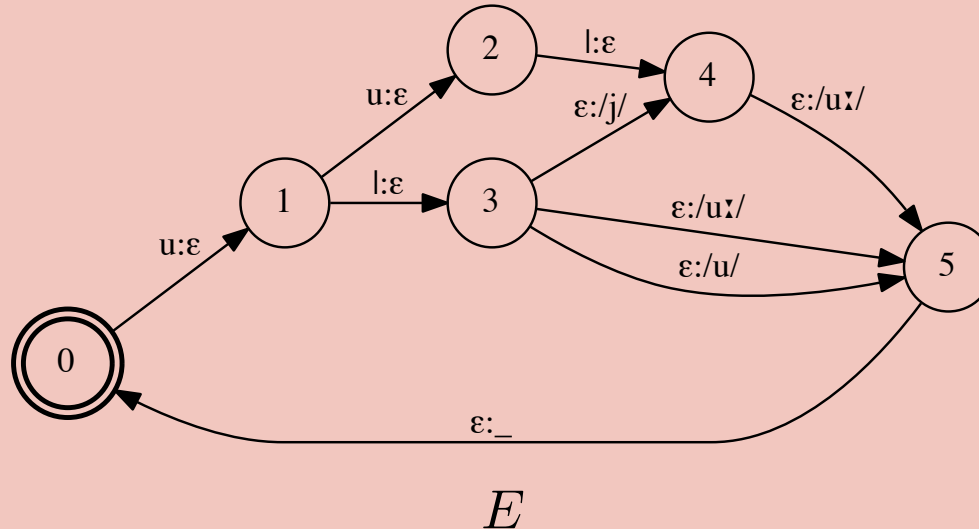
Alignment

- Construct letter FSTs W and T for a word w and its transcription t
- Alignment of w and t is generated by a series of compositions which filters out all non-matching pairings

$$A_{W,T} = \pi_2(W \circ I_G) \circ E \circ \pi_2(T \circ I_P)$$

Example

$$M = \{u:/u/, u:/u:/, u:/ju:/, uu:/u:/\}$$



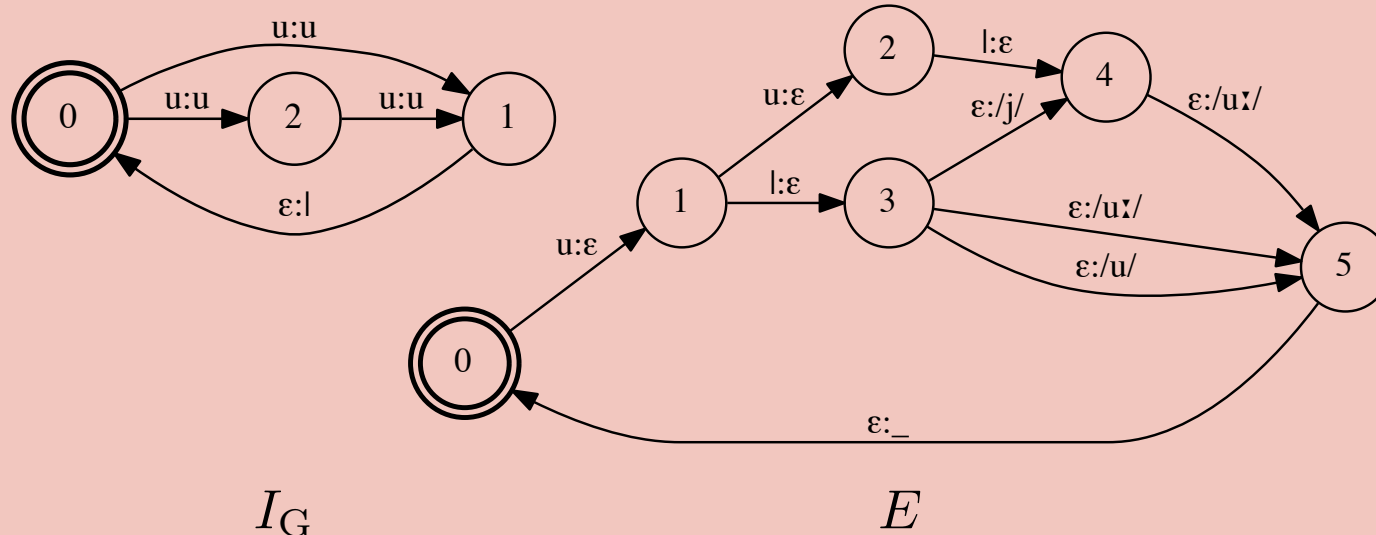
Alignment

- Construct letter FSTs W and T for a word w and its transcription t
- Alignment of w and t is generated by a series of compositions which filters out all non-matching pairings

$$A_{W,T} = \pi_2(W \circ I_G) \circ E \circ \pi_2(T \circ I_P)$$

Example

$$M = \{u:/u/, u:/u:/, u:/ju:/, uu:/u:/\}$$



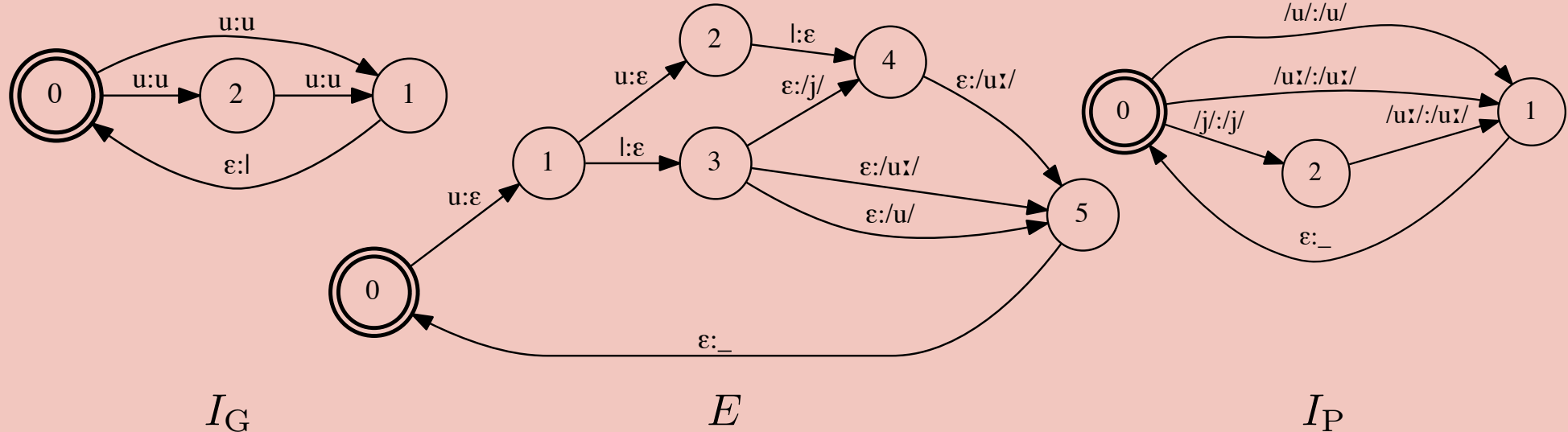
Alignment

- Construct letter FSTs W and T for a word w and its transcription t
- Alignment of w and t is generated by a series of compositions which filters out all non-matching pairings

$$A_{W,T} = \pi_2(W \circ I_G) \circ E \circ \pi_2(T \circ I_P)$$

Example

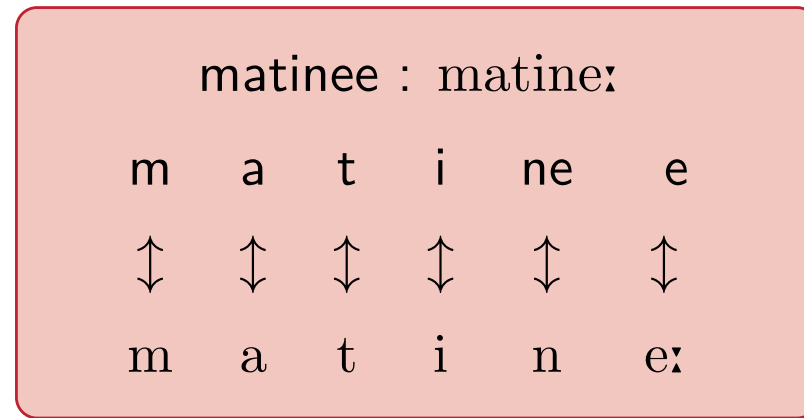
$$M = \{u:/u/, u:/u\:/, u:/ju\:/, uu:/u\:/\}$$



Alignment

Extended mappings

- Procedure allows for more complex mappings, i.e. context restriction
- Treatment of multiple alignments:



⇒ Conflicting rules may be disambiguated using lookahead conditions

Segmentation

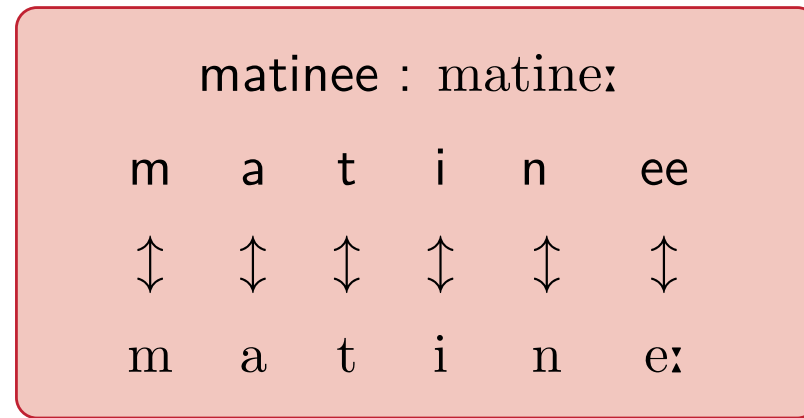
- I_G is used to generate possible grapheme level segmentations for subsequent transcription at runtime



Alignment

Extended mappings

- Procedure allows for more complex mappings, i.e. context restriction
- Treatment of multiple alignments:



⇒ Conflicting rules may be disambiguated using lookahead conditions

Segmentation

- I_G is used to generate possible grapheme level segmentations for subsequent transcription at runtime



Transcription

Idea

- Given aligned word-transcription pairs, transcription may be considered as *sequence labelling* problem
- Grapheme sequences are *observations*, phoneme sequences are *labels*
- Many existing methods, e.g. *Hidden Markov Models*, *Support Vector Machines*, *Conditional Random Fields* (cf. Erdogan 2010)

CRFs

(Lafferty et al. 2001)

- Graph-based model: labels and observations are represented by nodes
- Labelling is based on a set of random variables expressing characteristics of the observation \rightsquigarrow *features*
- Training process computes
 - Transition probabilities
 - Influence (weight) of the pre-defined features
- Runtime: Find the most likely state sequence



Transcription

Features

- Selection of features is a non-trivial task (i.e. no “inference” method)
- Given an input string $o = o_1 \dots o_n$, gramophone relays only on the (observable) grapheme context
 - Each position i is assigned a feature function f_j^k for each substring of o of length $m = (k - j + 1) \leq N$ within a context window of $N - 1$ characters relative to position i
 - N is the *context size window* or “order” of a gramophone model

$$f_j^k(o, i) = o_{i+j} \cdots o_{i+k} \text{ for } -N < j \leq k < N$$

$N = 1$	o_{i-3}	o_{i-2}	o_{i-1}	o_i	o_{i+1}	o_{i+2}	o_{i+3}
	m	a	t	i	n	e	e



Transcription

Features

- Selection of features is a non-trivial task (i.e. no “inference” method)
- Given an input string $o = o_1 \dots o_n$, gramophone relays only on the (observable) grapheme context
 - Each position i is assigned a feature function f_j^k for each substring of o of length $m = (k - j + 1) \leq N$ within a context window of $N - 1$ characters relative to position i
 - N is the *context size window* or “order” of a gramophone model

$$f_j^k(o, i) = o_{i+j} \cdots o_{i+k} \text{ for } -N < j \leq k < N$$

$N = 2$	o_{i-3}	o_{i-2}	o_{i-1}	o_i	o_{i+1}	o_{i+2}	o_{i+3}
	m	a	t	i	n	e	e



Transcription

Features

- Selection of features is a non-trivial task (i.e. no “inference” method)
- Given an input string $o = o_1 \dots o_n$, gramophone relies only on the (observable) grapheme context
 - Each position i is assigned a feature function f_j^k for each substring of o of length $m = (k - j + 1) \leq N$ within a context window of $N - 1$ characters relative to position i
 - N is the *context size window* or “order” of a gramophone model

$$f_j^k(o, i) = o_{i+j} \cdots o_{i+k} \text{ for } -N < j \leq k < N$$

$N = 3$	o_{i-3}	o_{i-2}	o_{i-1}	o_i	o_{i+1}	o_{i+2}	o_{i+3}
	m	a	t	i	n	e	e



Rating

Idea

- Select the “best” transcription from the segmented and labeled candidates
- Statistical model defined over strings of grapheme-phoneme segment pairs (“graphones”)
- N -gram model: joint probability as product of conditional probabilities under Markov assumptions

$$P(gp_0 \dots gp_n) \approx \prod_{i=0}^n P(gp_i | gp_{i-N} \dots gp_{i-1})$$

Implementation

- Interpolate all k -gram distributions with $1 \leq k \leq N$ (*Jelinek & Mercer 1980*)
- Combined with Kneser-Ney discounting for treatment of out-of-vocabulary items (*Kneser & Ney 1995*)
- Model parameters are estimated from (aligned) word-transcription pairs
- Implementable within the finite-state calculus (*Pereira & Riley 1997*)



Experiments

Corpora & Mappings

- **de-LexDB**: 71,481 words, 277 graphone types (Gibbon & Längen 2000)
- **de-Wiki** : 147,359 words, 589 graphone types (<http://de.wiktionary.org>)
- **en-CELEX**: 73,736 words, 463 graphone types (Baayen et al. 1995)

Method

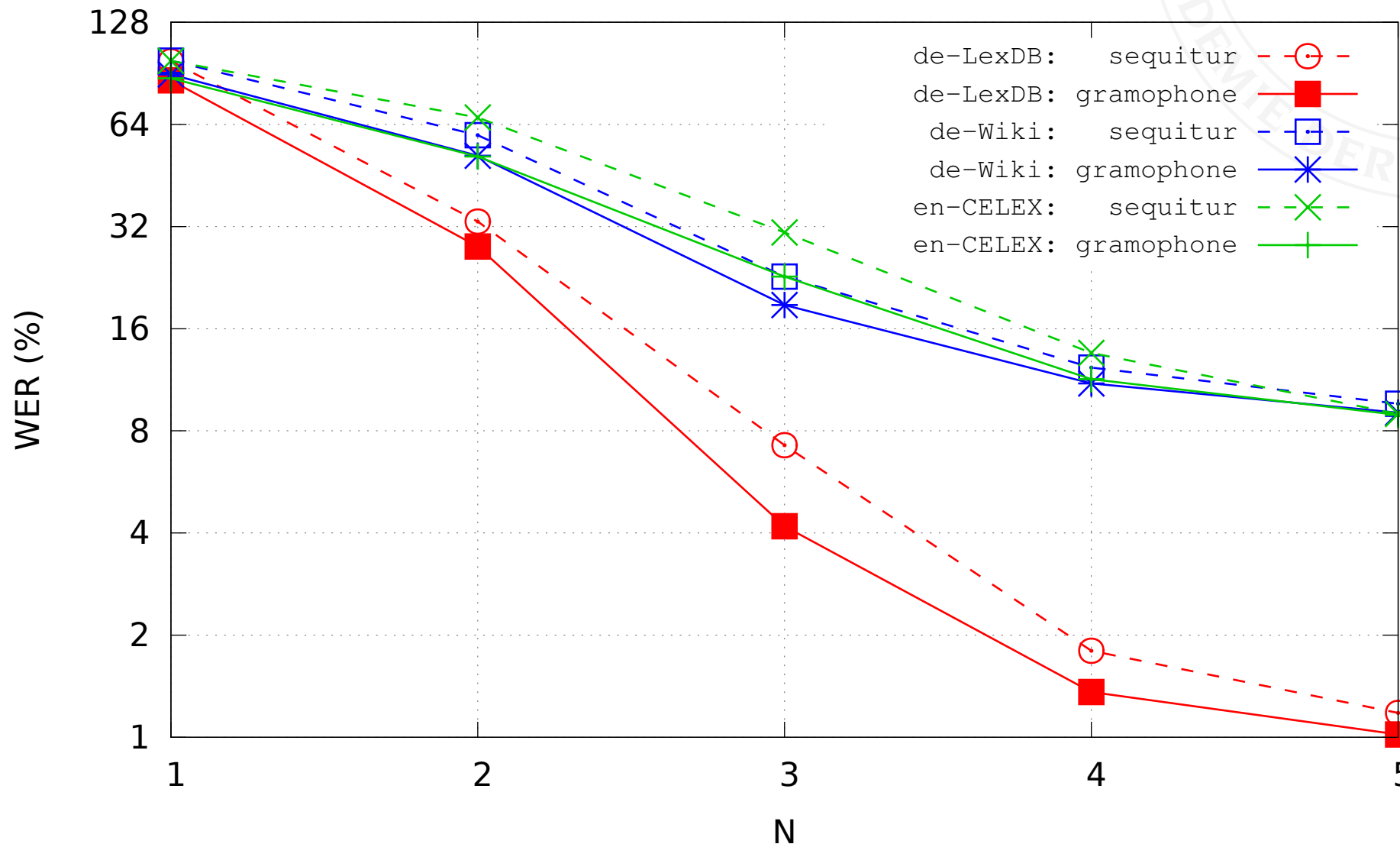
- Compare graphophone *versus* sequitur (Bisani & Ney 2008)
- Test model orders $N \in \{1, 2, 3, 4, 5\}$ using 10-fold cross validation
- Investigate both **word** and **phoneme** error rates (WER, PER)

Implementation

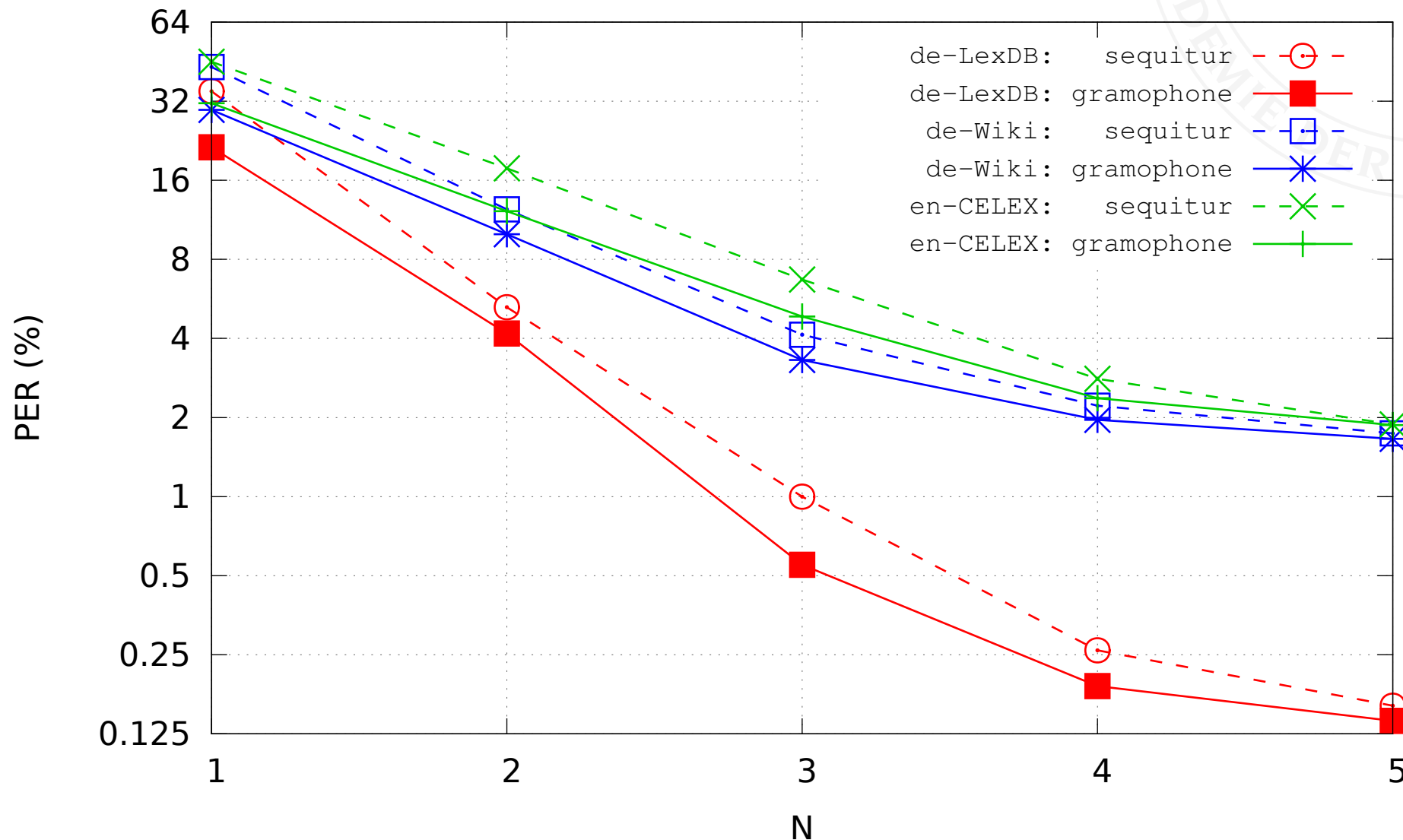
- OpenFST for alignment and segmentation (Allauzen et al. 2007)
- wapiti for CRF training and application (Lavergne et al. 2010)
- OpenGRM for candidate rating (Roark et al. 2012)



Results: Word Error Rate



Results: Phoneme Error Rate



Results: Discussion

General Trends

- gramophone outperformed sequitur for all conditions tested
- performance gain drops as model order increases, negligible for $N = 5$
- upper bound imposed by mapping heuristics beyond $N = 5$?
- LexDB performance looks suspiciously good
 - LexDB data were to a large extent automatically generated *(Längen p.c.)*

Interesting Phenomena

- de-Wiki: 25% of the phoneme errors concern schwa deletion
- de-Wiki: glottal stop is not a big issue
- en-CELEX: more uniform distribution of errors, largest class is schwa $\leftrightarrow V$ (22%)

	ən/n̩	əl/l̩	əm/m̩	ʔ/ Ɂ ʔ
seq	5114	756	307	172
gp	5010	633	299	146



Summary & Outlook

What We Did (instead of summer holidays)

- Novel conversion method based on three simple steps
 - Manually driven alignment/segmentation candidate generation
 - Candidate transcription with CRFs
 - Selection of the most likely candidate using N -gram LM
- Performance comparable to a state-of-the-art method

Still To Do

- Upper bound on performance imposed by segmentation heuristics (?)
- (Approximate) implementation using (weighted) finite-state methods (?)
 - Transducer (segmentation) \leftrightarrow pair acceptor (LM)
 - Linear chain CRFs \neq (W)FSTs
- Extensions
 - Integrate results of preceding morphological analysis
 - Predict syllabification, stress patterns





The End

/ðɪ ɛnd/

Thank you for listening!

