# Automatic classification of lexical stress errors for German CAPT

Anjana Vakil and Jürgen Trouvain

UNIVERSITÄT
DES
SAARLANDES

Department of Computational Linguistics and Phonetics
Saarland University, Saarbrücken, Germany

SLaTE 2015, Leipzig
4 September 2015

UNIVERSITÄT
DES
SAARLANDES

Lexical stress: Accentuation/prominence of syllable(s) in a word

In German:

- ► Variable placement, contrastive function

  um·FAHR·en     vs.     UM·fahr·en
  *to drive around*        *to run over*

- ► Reflected by duration, F0, intensity
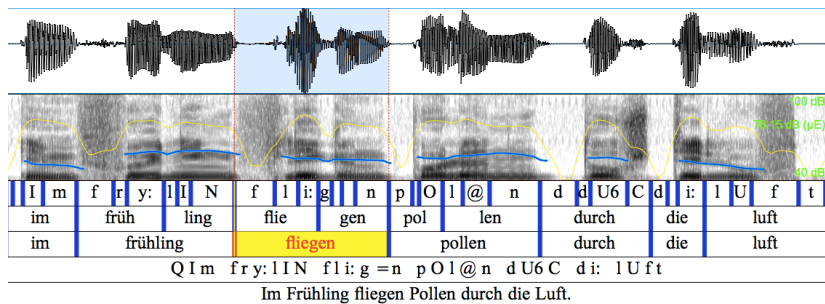- ► Impacts intelligibility of non-native (L2) speech

- Contrastive lexical stress (LS) difficult for French speakers
- CAPT can help; requires automatic diagnosis
- Classification of LS errors in L2 German unexplored

Classification of LS errors by French learners of German

*How feasible is it?*

*Which features are most useful?*

UNIVERSITÄT
DES
SAARLANDES

Subset of IFCASL corpus of French-German speech
(Fauth et al. 2014)



| I | m | f | r | y: | l | I | N | f | l | i: | g | n | p | O | l | @ | n | d | d | U6 | C | d | i: | l | U | f | t |
| im | | früh | | | ling | | | flie | | | gen | | pol | | | len | | | durch | | | | die | | luft | | |
| im | | frühling | | | | | | fliegen | | | | | pollen | | | | | | durch | | | | die | | luft | | |

Q I m   f r y: l I N   f l i: g =n   p O l @ n   d U6 C   d i:   l U f t

Im Frühling fliegen Pollen durch die Luft.

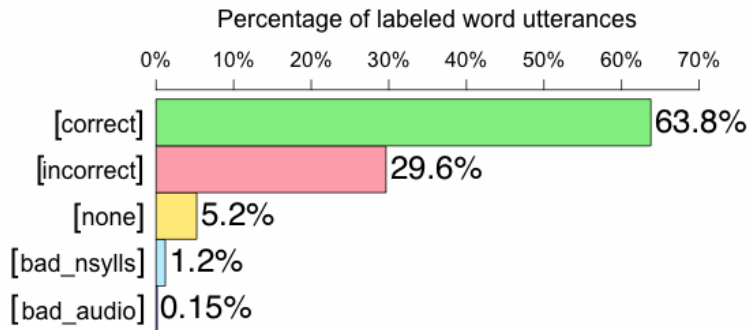Extracted utterances of 12 bisyllabic, initial-stress words

- ▶ 668 tokens from 56 French speakers - manually annotated
- ▶ 477 tokens from 40 German speakers - assumed correct

- Each token assigned a class label:

  [correct], [incorrect], [none]

  [bad_nsylls], [bad_audio]

- 15 annotators (12 native), each token labeled by $\geq 2$
- Varying phonetics/phonology expertise

Overall pairwise inter-annotator agreement

|  | Mean | Maximum | Median | Minimum |
|---|---|---|---|---|
| % Agreement | 54.92% | 83.93% | 55.36% | 23.21% |
| Cohen's $\kappa$ | 0.23 | 0.61 | 0.26 | -0.01 |

- ▶ Variability not explained by annotator L1 or expertise
- ▶ Single gold-standard label selected for each token

Train & evaluate CART classifiers using WEKA toolkit

## Training data
- Manually annotated L2 utterances
- Automatically annotated L1 utterances (all [correct])

## Held-out testing data
- Feature comparison: 1/10 of L2 utterances (random)
- Unseen speakers: all utterances from 1 of 56 L2 speakers

## Evaluation
- Compute agreement (% and $\kappa$) with gold standard
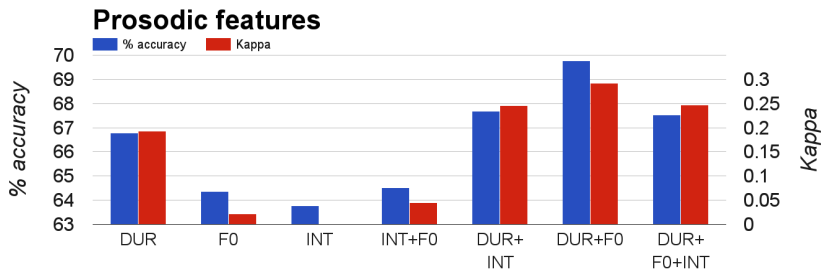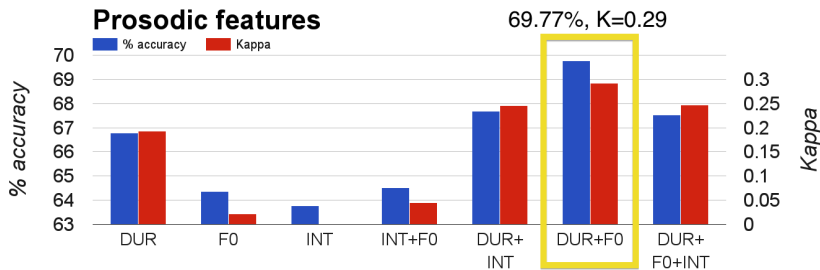- Cross-validation (10 or 56 folds)

# Feature sets

Prosodic feature sets

- ► DUR - Duration (relative syllable & nucleus lengths)
- ► F0 - Fundamental frequency (mean, max., min., range)
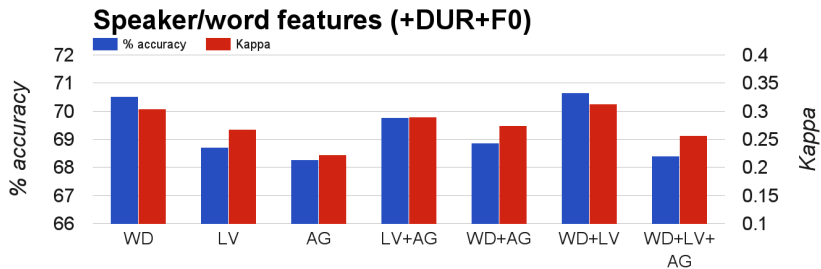- ► INT - Intensity (mean, max.)

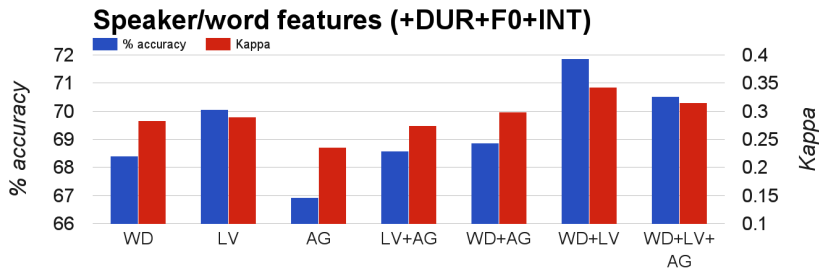Pitch and energy contours calculated using JSnoori software
(http://jsnoori.loria.fr)
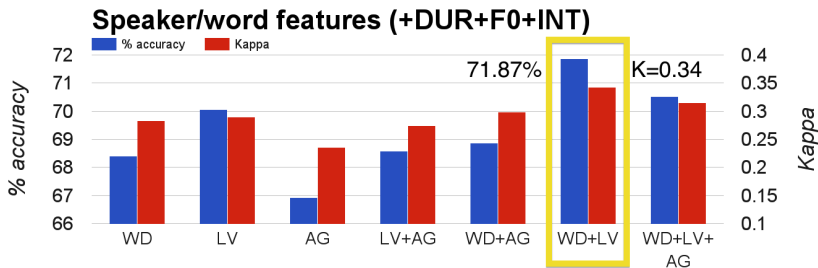
Other features

- WD - Word uttered (e.g. *Flagge*)
- LV - Speaker's CEFR skill level (A2|B1|B2|C1)
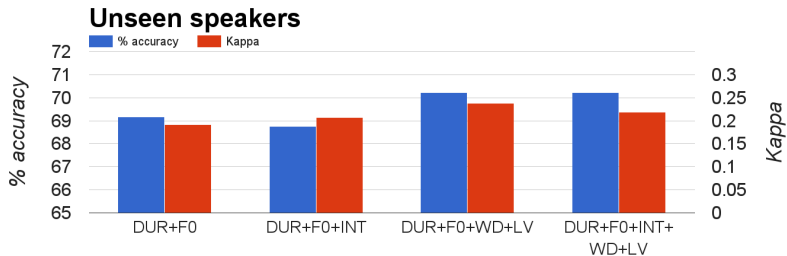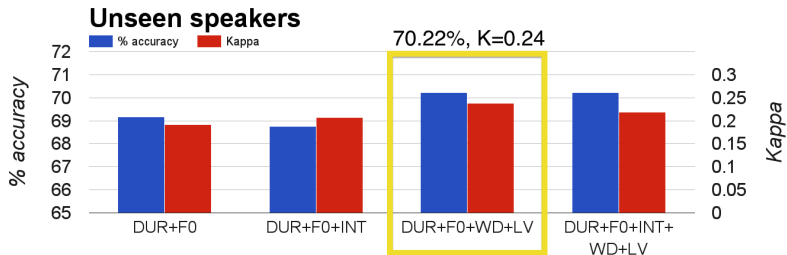- AG - Speaker's age/gender (Girl|Boy|Woman|Man)

**Prosodic features**

Speaker/word features (+DUR+F0)

Speaker/word features (+DUR+F0+INT)

Speaker/word features (+DUR+F0+INT)

UNIVERSITÄT
DES
SAARLANDES

|  | % agreement | $\kappa$ |
|---|---|---|
| Best classifier vs. gold standard | | |
|     Random test set | 71.87% | 0.34 |
|     Unseen speakers | 70.22% | 0.24 |
| Majority ([correct]) classifier vs. gold | 63.77% | 0.00 |
| Human vs. human | 54.92% | 0.23 |

- ▶ Results are encouraging in this context
- ▶ Still want better performance for real-world use

- Classification-based diagnosis of lexical stress errors
  novel approach in German CAPT
- Results of $>70\%$ accuracy encouraging
  (especially considering low human-human agreement)
- Still much room for improvement

Future directions
- More powerful machine learning algorithms
- Additional features (e.g. vowel quality, phrase information)
- Online, semi-supervised learning/active learning