

Automatic classification of lexical stress errors for German CAPT

Anjana Sofia Vakil

Department of Computational Linguistics & Phonetics

Saarland University, Saarbrücken, Germany

anjanav@coli.uni-saarland.de

Abstract

[TODO Abstract (200 word limit)]

[illegible]

Index Terms: computer-assisted pronunciation training, CAPT, word prosody, German [TODO are these OK?]

1. Introduction

For adult learners of a second language (L2), the phonological system of the L2 can pose a variety of difficulties. For certain L2s, such as German or English, one important difficulty involves the accurate prosodic realization of lexical stress, i.e. the accentuation of certain syllable(s) in a given word, with the placement of stress within a word varying freely and carrying a contrastive function in such languages [1]. Lexical stress is an important part of German word prosody, and has been found to have an impact on the intelligibility of non-native German speech [2]. Coping with this phenomenon in German is especially challenging for native (L1) French speakers, because lexical stress is realized very differently (or perhaps not at all) in the French language [3, 4].

To overcome this difficulty and improve their L2 word prosody, learners typically need to have their pronunciation errors pointed out and corrected by a language instructor; unfortunately, the lack of attention typically given to pronunciation in the foreign language classroom, along with other factors such as high student-to-teacher ratios, make this level of individualized attention not always feasible in a classroom setting [5, 6, 7]. Fortunately, advances in Computer-Assisted Pronunciation Training (CAPT) over recent decades have made it possible to automatically provide highly individualized analysis of

learners' prosodic errors, as well as feedback on how to correct them, and thus to help learners achieve more intelligible pronunciation in the target language. However, while much research has gone into the creation and improvement of CAPT systems for English (see e.g. [8, 9]), relatively little work has been done on the development of CAPT systems for German, especially on those targeting errors in German prosody.

This paper describes work that advances the state of German CAPT by applying machine learning methods to the task of diagnosing lexical stress errors in non-native German speech, a necessary prerequisite for delivering individualized corrective feedback on such errors in a CAPT system. The paper is organized as follows: Section 2 provides background on the phenomenon of lexical stress as it is realized in German and French word prosody, motivates the creation of CAPT systems that address this error specifically, and summarizes some past work related to this topic. Section 3 describes the manual annotation of lexical stress errors in a small corpus of L2 German speech, carried out to create labeled training and test data for the classification experiments explained in section 4. Section 5 presents and analyzes the results of these experiments. Finally, section 6 offers some concluding remarks and outlines possible directions for future work.

2. Background and related work

Broadly speaking, lexical stress is the phenomenon of how a given syllable is accentuated within a word [1], i.e. how a syllable is given a more prominent role such that this syllable is perceived as “standing out” [10]. This perceived prominence of a syllable is a function not merely of the segmental characteristics of the uttered syllable, i.e. the speech sounds it contains, but rather of its (relative) suprasegmental properties, namely:

- duration, which equates on the perceptual level to length;
- fundamental frequency (F0), which corresponds to perceived pitch; and
- intensity (energy or amplitude), which perceptually equates to loudness.

In variable-stress languages, such as German and English, the location of lexical stress in a word is not always predictable, and therefore knowing a word requires, in part, knowing its stress pattern. This allows lexical stress to serve a contrastive function in these languages, e.g. distinguishing *UMfahren* (to drive around) from *umFAHren* (to run over with a car) in German. Furthermore, in German, misplaced stress can disrupt understanding even in cases where there is no stress-based minimal pair [2]. However, in fixed-stress languages, stress is completely predictable, as it always falls on a certain position in the word (e.g. the final syllable), making the lexical stress pattern less crucial to the knowledge of a word than in variable-stress

languages. Furthermore, in fixed-stress languages there may be a weaker distinction between stressed and unstressed syllables. While French has often been categorized as a fixed-stress language, given that word-final syllables are given prominence when a French word is pronounced in isolation, some argue that it may be more properly considered a language without lexical stress, in that speakers do not seem to accentuate any syllable within the word, with word-final lengthening effects explained by interactions with the realization of phrasal accent (lengthening of the final syllable in each prosodic group or phrase) [3, 4]. Regardless, French has no contrastive word-level stress [4, p. 89], and in this respect differs considerably from German.

This difference between the languages leads us to expect French learners of German to have difficulties with both perception and production of lexical stress prosody. Although little research has been done on the nature of lexical stress errors for this particular L1-L2 pair, Hirschfeld and Trouvain [7] report that such errors are commonly observed in German spoken by French natives. Research on French speakers' perception of Spanish, another contrastive-stress language, has revealed that these speakers seem to be "deaf" to lexical stress, i.e. seem to have significant and lasting difficulty perceiving and remembering stress contrasts [3]. With respect to production, studies of L2 Dutch have shown that French speakers, especially beginners, make systematic errors with lexical stress, exhibiting a tendency to stress the final syllable of Dutch words even when stress should be placed on the initial or medial syllable [11, 4]. Similar findings have also been reported for French learners of English [12]. The high (anticipated) frequency of lexical stress errors in the speech of this L1-L2 group is thus one motivating factor for the creation of CAPT systems to help learners identify and correct such errors.

Another motivation behind this work's focus on lexical stress errors is the high impact such errors may have on the intelligibility of L2 German speech. Intelligibility, as opposed to lack of a foreign accent, is generally considered to be the most important goal of pronunciation training [13, 5, 6, 14, 9]. The exact definition of *intelligibility* is a topic of debate, but here we will follow Munro and Derwing [13, p. 289] in understanding it broadly as "the extent to which a speaker's message is actually understood by a listener." Generally speaking, prosodic errors have often been found to have a larger impact on the perceived intelligibility of L2 speakers than segmental errors (Derwing and Munro, 2005; Hahn, 2004; Witt, 2012), and several studies have found lexical stress errors to have a particularly strong impact on intelligibility in free-stress languages like English and Dutch [1, 14]. Though relatively little research has been done on how various pronunciation errors affect intelligibility in L2 German specifically, some studies suggest that lexical stress errors may hinder intelligibility of L2 German speech more than other types of errors [2, 7]. Stress errors may also affect perception of segmental errors in the L2 learners' speech; for example, segmental errors occurring in stressed syllables may be more noticeable than those in unstressed syllables [1, 11]. It would therefore seem that there may be a strong connection between lexical stress errors and intelligibility in L2 German speech, though more research is needed to clarify the nature of this relationship. **[TODO remove that sentence?]**

Though the frequency and impact of lexical stress errors in the speech of French learners of German thus constitute strong reasons to develop CAPT tools to treat such errors, in order for such systems to be viable, the feasibility of reliable automatic detection of this type of error must be demonstrated. **[TODO Make rest of this par about how comparison-based**

diagnosis is usually used, then start new par with following sentence?] To our knowledge, no work has been reported on automatic classification-based diagnosis of lexical stress errors in L2 German speech, but in recent years machine learning methods have been applied with apparent success to the classification of lexical stress patterns in English words. Kim and Beutnagel [15] experimented with various classifiers to identify stress patterns in high-quality recordings of 3- and 4-syllable English words, reporting accuracy in the 80-90% range; in pilot experiments with low-quality recordings, however, the authors report lower accuracy: 70-80% on L1 speech and 50-60% on utterances by L2 speakers. Similarly, Shahin et al. [16] trained Neural Networks to classify stress patterns in bisyllabic words uttered by L1 English children, and reported classification accuracy over 90% for some stress patterns; though this work was conducted with a view to treating childhood L1 dysprosody, its relevance to our intended application of L2 CAPT is nonetheless clear. Building on these related investigations, this paper seeks to further explore the viability of automatic classification-based detection of lexical stress errors, with a particular focus on those made by French speakers of German. To this end, a small corpus of learner speech was manually annotated for lexical stress errors, as described in section 3. Using the resulting labeled L2 data, in addition to data from L1 German speakers, a series of supervised machine learners were trained using a variety of representations of the prosodic and other features of each word utterance (see section 4.1, and these classifiers were evaluated with reference to the manually-produced labels of held-out test data (see section 4.2). Section 5 presents and analyzes the results of these evaluations.

3. Data

Error-annotated speech data from German learners is a prerequisite for the supervised training and evaluation of classifiers for lexical stress realizations in L2 German speech, yet to our knowledge no corpus of learner German with such annotation is publicly available. To fill this need, as well as to shed light on the perception of lexical stress errors in L2 German speech, a small corpus of speech by L1 French learners of German was manually annotated for such errors by native and non-native German speakers with varying levels of phonetics/phonology expertise. This section describes the data selected for annotation (section 3.1) and the method by which lexical stress realizations in this data were annotated (section 3.2), and presents an analysis of the observed inter-annotator agreement (section 3.3) and distribution of errors (section 3.4) in the annotated dataset.

3.1. The IFCASL corpus of learner speech

The learner speech data used in this work has been excerpted from the IFCASL corpus [17, 18], a collection of phonetically diverse utterances in French and German spoken by both native speakers and non-native speakers with the other language as L1. This is the first known corpus of L2 speech in both directions of the French-German language pair, and is thus an invaluable resource for research on pronunciation errors **[TODO between these languages]**.

The IFCASL corpus contains recordings of approximately 50 L1 speakers of each language reading carefully constructed sentences (and a short text) in both languages, such that both L1 and L2 speech was recorded for each speaker. Each L1 speaker group has an even gender distribution, and contains approximately 10 children (adolescents of 15-16 years of age) and 40

Table 1: Word types selected from the IFCASL corpus for lexical stress error annotation. Canonical pronunciations for each word type are given in IPA notation. The rightmost column lists the number of tokens (utterances) of each word type in the annotated dataset.

Word	Pronunciation	Part of speech	English meaning	Tokens
E-mail	/ˈiː.meɪl/	noun	e-mail	56
Flagge	/ˈfla.ɡə/	noun	flag	55
fliegen	/ˈfliː.ɡn/	verb	to fly	56
Frhling	/ˈfryː.lɪŋ/	noun	spring (season)	56
halten	/ˈhal.tn/	verb	to hold	56
manche	/ˈman.ʃə/	pronoun	some	56
Mrder	/ˈmɔ̃r.dɐ/	noun	murderer	56
Pollen	/ˈpɔ̃.lən/	noun	pollen	55
Ring	/ˈʁɪŋ.ən/	noun	rings	55
Tatort	/ˈtʰat.ʔɔ̃t/	noun	crime scene	56
tragen	/ˈtʰrʰaː.ɡn/	verb	to wear	55
Tschechen	/ˈtʃɛ.ʃn/	noun	Czechs	56

adults. A variety of self-reported L2 proficiency levels are also represented in the corpus: the recorded adults span levels A2 (beginner) through C1 (advanced) of the Common European Framework of Reference [\[TODO footnote url\]](#), the children levels A2 (beginner) and B1 (low intermediate).

While L2 French speech is thus also captured in the IFCASL corpus, the annotation effort described here focuses exclusively on the German-language subset of the corpus. Only utterances from the sub-corpus of L2 German speech by L1 French speakers (henceforth IFCASL-FG) were manually annotated; native utterances from the L1 German sub-corpus (IFCASL-GG) were assumed to contain only correct lexical stress realizations.

The subset of IFCASL-FG selected for manual error annotation, (henceforth simply the dataset) consists of utterances of twelve bisyllabic word types (see table 1), each of which has primary stress on the initial syllable. Only bisyllabic words were selected to simplify comparison between stressed and unstressed syllables, and only initial-stress words because this is the stress pattern which native (L1) French speakers are expected to have the most difficulty producing in German, given the phenomenon of final lengthening in French (see section 2).

In addition to the recordings themselves, the IFCASL corpus contains phone- and word-level segmentations of each utterance, produced automatically by forced alignment with the corresponding text prompts [18]. Although the corpus also contains manual corrections of these segmentations, the work reported here relies exclusively on the automatically-generated segmentations to more accurately represent the conditions of a real-world CAPT system, which would not have recourse to manual verification of the phone or word boundaries identified by the aligner. Using these segmentations, tokens (utterances) of each selected word type were extracted from the recorded segmentations automatically using Praat[\[TODO footnote url\]](#); counts of the available tokens for each word type are listed in table 1. The dataset annotated for lexical stress errors comprises 668 word tokens in total. Five tokens had to be excluded from the dataset, as sentence-level disfluencies (e.g. false starts

or repetitions of phrases) prevented accurate automatic extraction of the word utterance; a fully-fledged CAPT system would need to deal with such disfluencies automatically, e.g. with a pre-processing step which detects disfluencies and prompts the learner to re-record their utterance if needed (see e.g. [19]).

3.2. Annotation method

The annotation task consisted of assigning one of the following labels to each word token (utterance) in the dataset described in the previous section:

- [correct]: the speaker clearly stressed the correct (initial) syllable
- [incorrect]: the speaker clearly stressed the incorrect (final) syllable
- [none]: the speaker did not clearly stress either syllable, or the annotator was unable to determine which syllable was stressed
- [bad_nsylls]: the speaker pronounced an incorrect number of syllables (e.g. inserted an extra syllable), making it impossible to judge if stress was realized correctly
- [bad_audio]: a problem with the audio file (e.g. noise or inaccurate segmentation) interfered with the annotator’s ability to judge the stress realization

Annotation was performed using a graphical tool scripted in Praat. This tool displayed the given word’s text, and allowed the annotator to listen to the given word utterance and the sentence utterance from which it was extracted as many times as they wished. Once they had reached a judgment about the lexical stress realization of the utterance, the annotator had to click one of five buttons, corresponding to the possible labels, to record their judgment. A single annotation session consisted of annotating all 55-56 tokens of each of three word types, and lasted approximately 15 minutes.

A total of 15 annotators participated, varying with respect to their L1 and level of phonetics/phonology expertise. The native languages represented included German (12 annotators), English (2), and Hebrew (1); the L1 English and Hebrew speakers all speak German as [\[TODO a\]](#) L2. In terms of expertise, the annotators were broadly categorized as *experts* (professional phonetics/phonology researchers), *intermediates* (university students enrolled in an experimental phonology course), or *novices* (those with negligible phonetics/phonology training or experience annotating speech data). Among the 15 annotators, there were two experts, 10 intermediates, and three novices.

Each annotator was assigned three word types to annotate in a single session, with the exception of one who annotated six word types over two sessions. Assignments ensured that each word token was annotated by at least two native German speakers, and to maximize the amount of overlap between annotators in order to obtain as many pairwise measures of annotator agreement as possible (see section 3.3).

3.3. Inter-annotator agreement

Any evaluation of an automatic error detection system, including that described in this work, should be performed with an understanding of the difficulty of the error-detection task for human listeners. To obtain a clearer picture of this task, we therefore conducted an analysis of the inter-annotator agreement observed in the annotations collected as described in the previous section. If human annotators often disagree about whether a

Table 2: Overall pairwise agreement between annotators

	% Agreement	Cohen’s κ
Mean	54.92%	0.23
Maximum	83.93%	0.61
Median	55.36%	0.26
Minimum	23.21%	-0.01

given L2 utterance contains a lexical stress error, this may indicate that the task is a difficult one, thus encouraging a more lenient evaluation of an automatic error-detection system. However, if humans are generally in strong agreement, this may reflect a lower level of difficulty, and give reason to judge the performance of an automatic system by a higher standard.

For 268 of the 668 utterances annotated, i.e. approximately 40% of the dataset, annotators were unanimous in their label assignments; for the other 400 utterances (60%), at least one annotator chose a different label than the other(s) who annotated the same utterance. To make sense of these differences, agreement in label assignments was calculated for each pair of annotators who overlapped, i.e. labeled any of the same tokens. Pairwise agreement was quantified in terms of percentage agreement (i.e. the number of tokens to which the two annotators assigned the same label, divided by the total number of tokens they both annotated), and Cohen’s Kappa (κ) statistic [20]. To obtain an overall measure of inter-annotator agreement for the entire annotated dataset, the agreement between each pair of overlapping annotators was calculated, and the minimum, median, mean, and maximum values over all pairwise comparisons were computed; these values are given in table 2.

This simple analysis reveals a few interesting observations. First, the mean and median percentage agreement values near 55% indicate that annotators seem to agree about the accuracy of lexical stress realizations just slightly more than they disagree, and the mean and median κ values near 0.25 characterizes the overall agreement as “fair” in the Landis and Koch schema [21]. However, the minimum and maximum κ values reveal that agreement between different pairs of annotators ranges from “poor” to “substantial” [21], as also reflected in the correspondingly large gap between the minimum and maximum percentage agreement observed. On the whole, then, it appears that inter-annotator agreement in this error annotation task is relatively low, though there seems to be considerable variation between individual annotators. This may simply signal that (some of) the particular annotators participating in this study are not very reliable in their judgments of lexical stress accuracy, but it may also indicate that diagnosing errors in L1 French speakers’ realizations of lexical stress in German is a difficult task, even for humans.

3.4. Error distribution

Though the comparison of the (potentially different) labels assigned to each word utterance by different annotators is interesting, a single “gold-standard” label for each utterance ultimately had to be determined, as a representation of the ground truth with which to train and evaluate the automatic error classifier(s) [TODO rewrite that sentence]. In some cases, assigning a gold-standard label was trivial, e.g. when all or a majority of annotators agreed. However, in other cases a choice had to be made between competing candidate labels. [TODO The decision-making procedure] prioritized experts’ judgments,

Table 3: Overall frequency of lexical stress errors in the annotated data

Label	Tokens	% of corpus
correct	426	63.77%
incorrect	198	29.64%
none	35	5.24%
bad_nsylls	8	1.20%
bad_audio	1	0.15%
Total	668	100%

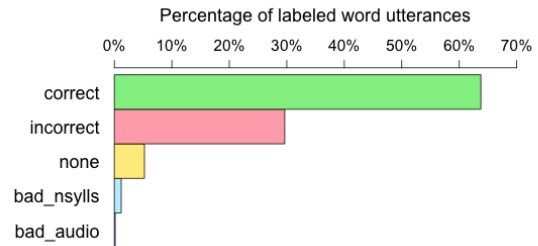


Figure 1: Overall distribution of lexical stress errors in the annotated data

avored certain judgments ([correct] or [incorrect]) over uncertain ([none]), and gave learners the benefit of the doubt when annotators disagreed as to whether the utterance was [correct] or [incorrect].

To shed further light on the error classification task, the overall distribution of lexical stress errors in the annotated dataset was analyzed with reference to the gold-standard labels thus determined. As seen in table 3 and illustrated in fig. 1, most (63.77%) of learners’ lexical stress productions were judged to be correct; in other words, almost two-thirds of the time, learners clearly stressed the correct (initial) syllable in the uttered word. However, learners also seemed to make mistakes regularly, with 29.64% of their word utterances labeled [incorrect] and another 5.24% labeled [none]. If we consider both [incorrect] and [none] utterances as types of lexical stress errors, then errors were observed in just over one-third of utterances. This considerable proportion of errors seems to confirm the expectation (mentioned in section 2) that French learners of German frequently make lexical stress errors.

4. Method

As mentioned in section 2, the classification-based approach to identifying lexical stress errors has not been sufficiently explored in CAPT research, especially research on CAPT for German. By way of a preliminary investigation of the feasibility of this type of error diagnosis, a series of classification experiments were conducted in an effort to determine:

- how well lexical errors can be identified by a classification-based approach, in comparison to the accuracy of human listeners in identifying such errors,
- [TODO which of the prosodic features discussed in ??] are most useful for classification, and
- whether a classification-based approach can lead to reasonably accurate diagnosis for words or speakers not

seen in the training data. [TODO explain]

The WEKA machine learning toolkit[TODO footnote URL] was used to train and evaluate classifiers for these experiments. In this work, only simple Classification And Regression Tree (CART) classifiers were used. Many other classification algorithms are implemented in WEKA, some of which could conceivably offer better performance; here, CARTs were chosen for their simple training process and their ease of interpretation by humans. In future work (see section 6), it would be interesting to compare different classification algorithms to see if other classifiers are more effective for this type of data, along the lines of the experiments by Kim and Beutnagel [15].

[TODO paraphrase the below]

For each relevant configuration (see sections 4.1 and 4.2), a CART is trained to classify utterances as belonging to one of the five label categories described in section 3.2. In practice, however, these trees classify every utterance as either [correct] or [incorrect], neglecting [none] and the other labels due to their comparatively low frequency in the data. Overall classification accuracy on the annotated sub-corpus was assessed by using held-out portions of the annotated data as test sets, and performing cross-evaluation on multiple train/test splits of the data. The features and divisions of data used in each experiment are described in the sections below.

4.1. Feature sets

4.2. Datasets for training and testing

5. Results

5.1. Feature performance

5.2. Performance on unknown words

5.3. Performance on unknown speakers

6. Conclusions and future work

7. References

- [1] A. Cutler, "Lexical Stress," in *The Handbook of Speech Perception*, D. B. Pisoni and R. E. Remez, Eds., 2005, pp. 264–289.
- [2] U. Hirschfeld, *Untersuchungen zur phonetischen Verständlichkeit Deutschlernender*, ser. Forum Phonetikum, 1994, vol. 57.
- [3] E. Dupoux, N. Sebastián-Gallés, E. Navarette, and S. Peperkamp, "Persistent stress 'deafness': The case of French learners of Spanish," *Cognition*, vol. 106, pp. 682–706, 2008.
- [4] M.-C. Michaux and J. Caspers, "The production of Dutch word stress by Francophone learners," in *Proceedings of the Prosody-Discourse Interface Conference 2013 (IDP-2013)*, 2013, pp. 89–94.
- [5] A. Neri, C. Cucchiari, H. Strik, and L. Boves, "The pedagogy-technology interface in computer assisted pronunciation training," *Computer Assisted Language Learning*, 2002.
- [6] T. M. Derwing and M. J. Munro, "Second Language Accent and Pronunciation Teaching: A Research-Based Approach," *TESOL Quarterly*, vol. 39, no. 3, pp. 379–397, 2005.
- [7] U. Hirschfeld and J. Trouvain, "Teaching prosody in German as foreign language," in *Non-Native Prosody: Phonetic Description and Teaching Practice*, J. Trouvain and U. Gut, Eds. Walter de Gruyter, 2007, pp. 171–187.
- [8] M. Eskenazi, "An overview of spoken language technology for education," *Speech Communication*, vol. 51, no. 10, pp. 832–844, Oct. 2009.
- [9] S. M. Witt, "Automatic error detection in pronunciation training: Where we are and where we need to go," in *Proceedings of the International Symposium on Automatic Detection of Errors in Pronunciation Training (IS ADEPT)*, 2012, pp. 1–8.
- [10] G. Dogil and B. Williams, "The phonetic manifestation of word stress," in *Word Prosodic Systems in the Languages of Europe*, H. van der Hulst, Ed. Berlin: Walter de Gruyter, 1999, ch. 5, pp. 273–334.
- [11] M.-C. Michaux, "Exploring the production and perception of word stress by French-speaking learners of Dutch," in *Workshop on Crosslinguistic Influence in Non-Native Language Acquisition*, 2012.
- [12] A. Bonneau and V. Colotte, "Automatic Feedback for L2 Prosody Learning," in *Speech and Language Technologies*, I. Ipsic, Ed. InTech, 2011, no. 1977.
- [13] M. J. Munro and T. M. Derwing, "Foreign accent, comprehensibility, and intelligibility in the speech of second language learners," *Language Learning*, vol. 49, no. Supplement s1, pp. 285–310, 1999.
- [14] J. Field, "Intelligibility and the Listener: The Role of Lexical Stress," *TESOL Quarterly*, vol. 39, no. 3, p. 399, Sep. 2005.
- [15] Y.-J. Kim and M. C. Beutnagel, "Automatic assessment of American English lexical stress using machine learning algorithms," in *SLaTE*, 2011, pp. 93–96.
- [16] M. Shahin, B. Ahmed, and K. Ballard, "A neural network based lexical stress pattern classifier," *Qatar Foundation Annual Research Forum Proceedings*, no. 2012, p. CSP22, Oct. 2012.
- [17] J. Trouvain, Y. Laprie, B. Möbius, B. Andreeva, A. Bonneau, V. Colotte, C. Fauth, D. Fohr, D. Jouvét, O. Mella, J. Jügler, and F. Zimmerer, "Designing a bilingual speech corpus for French and German language learners," in *Corpus et Outils en Linguistique, Langues et Parole: Statuts, Usages et Méusages*, no. ii, Strasbourg, France, 2013, pp. 32–34.
- [18] C. Fauth, A. Bonneau, F. Zimmerer, J. Trouvain, B. Andreeva, V. Colotte, D. Fohr, D. Jouvét, J. Jügler, Y. Laprie, O. Mella, and B. Möbius, "Designing a Bilingual Speech Corpus for French and German Language Learners: a Two-Step Process," in *9th Language Resources and Evaluation Conference (LREC)*, Reykjavik, Iceland, 2014, pp. 1477–1482.
- [19] L. Orosanu, D. Jouvét, D. Fohr, I. Illina, and A. Bonneau, "Combining criteria for the detection of incorrect entries of non-native speech in the context of foreign language learning," in *SLT 2012 - 4th IEEE Workshop on Spoken Language Technology*, Dec. 2012.
- [20] J. Cohen, "A Coefficient of Agreement for Nominal Scales," *Educational and Psychological Measurement*, vol. 20, no. 1, pp. 37–46, Apr. 1960.
- [21] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33, no. 1, pp. 159–174, 1977.