

AN INVESTIGATION OF THE PRONUNCIATION OF ENGLISH WORDS IN GERMAN SMS TEXTS

Hongwei Ding and Oliver Jokisch

Institut für Akustik und Sprachkommunikation

Hongwei.Ding@ias.et.tu-dresden.de

Kurzfassung: This paper reports on pronunciation variations by German speakers of English words in German SMS speech database. The text database contains SMS communication in different domains, such as on holidays, families, congratulations, schools, etc. As expected, there exist many English words in the German texts. The speech database comprises 100 German speakers, these speakers are chosen with balanced age, female-male difference and dialect region. There are 13 different promptsheets for the speakers to read, so that only 7 or 8 speakers share the same promptsheet. Due to various characteristics of the English words and different English levels of the speakers, the English words in German SMS texts are pronounced in different ways, from standard pronunciation with an English accent to very poor mispronunciation with a very strong German accent. It is thus important to find out possible pronunciation variations of each word by average German speakers. These pronunciation variants can be built into the lexicon for the speech recognition of English words with German accent. It is further interesting to make a statistical study on the realization of particular English phonemes by average German speakers, the information can be exploited to derive pronunciation rules that can be applied to new vocabulary. We have investigated the pronunciation variations of all words that can have potential English pronunciations both at word level and phoneme level, only selected examples will be presented in this paper.

1 Introduction

Because of the worldwide communication, globalization of national economies, more and more English words can be found in German language, they appear in advertisements, such as "Douglas: Come in and find out" [7]; they appear in technical field such as "software"; they appear in teenage speech such as "cool"; they appear in news reports on economics such as "shareholder", on sport such as "trainer", on music titles such as "Yesterday once more", etc. And they are further reflected in the everyday communication, such as SMS texts. The morphological and grammatical changes of these English words in German have been widely studied in many researches [5, 3]. However, the pronunciation of these English words have received attention only in recent years [1].

With the development of speech technology, the pronunciation of English words becomes interesting for experts in phonetics and speech technologies. Some English words have been nativised for a long time, and will be pronounced according to German rules, such as "Trainer"; Other English words are regarded as original English, and English pronunciation can be expected, such as musical album "In My Mind". Sometimes because of the German context speakers try to pronounce the inserted English words in such way that they can be integrated into the German contexts, and sometimes because of unprofessional English level of the speakers English words will be pronounced according to German pronunciation rules. Anyway, English

words will be articulated in original English pronunciation, or nativised in German pronunciation, or somewhere between them. This paper investigates the pronunciation of 100 German native speakers of the English words appearing in German SMS texts, the statistics will be made at word level and phoneme level to provide an overview on the pronunciation variations, and a discussion on the investigation will also be provided.

An important application of the statistical results can be found in speech technology. The statistics can provide information for speech recognition on different pronunciation variants of English words with German accent, and the results can also provide implications for speech synthesis to deal with English words in German texts.

2 Linguistic Background

It is important to understand different kinds of loan words before we begin with the investigation.

2.1 Loan Words

Due to speech contact a great number of foreign words have been borrowed in German language in one way or the other. The following figure (Figure 1) from Duckworth [4] can illustrate different methods of loan words.

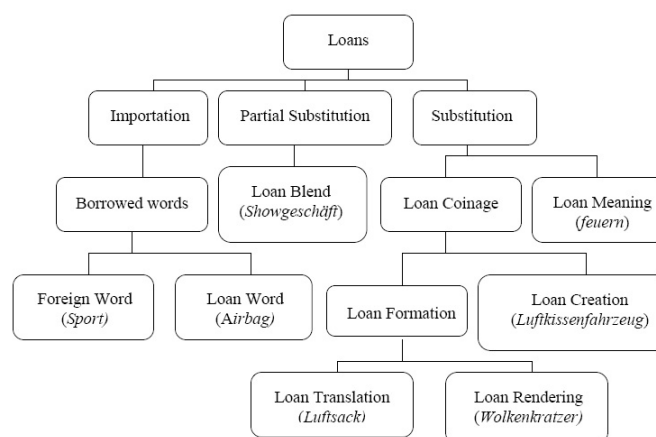


Figure 1 - Loan word categorization according to Duckworth [4]

For the investigation in this paper, we have only interest in those words that contain original English words, which can be potentially pronounced in English. In Figure 1, they include the types of "importation" and "partial substitution". Therefore English loan words which contain part or all English origins have been selected, and the pronunciation of different speakers would be transcribed and statistics would be made.

2.2 Comparison of English and German Phoneme Systems

Comparing English and German phoneme systems, we find they have some phonemes in common, our interest is directed to those English phonemes which do not exist in German phoneme system.

Some English consonants such as "/w/, /T/, /D/, /dZ/" can not be found in German phoneme system. Other consonants exist in German phoneme system, but have different distribution,

for example, "tS" appears not at syllable initial but only at syllable final. Because of "Auslautverhärtung" only devoiced obstruent can appear at syllable-final in German. The differences can be listed in great detail, but we do not have to exhaust all the differences, we only have to obtain some orientations for our analysis.

In SAMPA system some English and German vowels share the same symbols, such as "/u:/, /U/, /i:/, /I/", but they still have a little articulatory differences in aspects of tongue position and lip rounding, such substitution will generally be neglected in this investigation. Apart from that, there are some vowels, such as "/A:/, /Q/" and diphthongs "/eI/, /@U/", which are usually replaced with similar German vowels, such replacement will be considered in this investigation.

Since British pronunciation is widely accepted as standard in Germany, British lexicon is selected as reference for transcription and pronunciation.

3 Data Description

We recorded actually 100 speakers, but they were divided into 13 groups, there were about 7 or 8 persons in one group, different groups were asked to read different SMS texts. That means, the same texts would normally be read by 7 or 8 different speakers. The speakers were quite balanced in age, dialect region and female-male difference. Young speakers were normally students whose English levels were generally good. Older readers preferred nativised German pronunciation to original English pronunciation.

3.1 Selected English Words

The selected English loan words are listed in the following table (Table 1):

Some of the loan words belong to "importation", such as "airport, travel and work", others are "partial substitution", such as "Snowboardkurs, Cowboystiefeln". These loan words exist actually in everyday communication, the speakers should read them in the way they usually do in their speech.

As expected most of them are names of place (e.g. Southampton, Colorado), names of person (e.g. Harry Potter, James Bond), or names of music title (e.g. Little Miss Sunshine). The annotation of the speech database was decided according to the perception by a phonetic expert. The phone set for annotation includes all English phonemes and German phonemes. In the course of annotation, besides waveform reference, spectrogram reading, sound listening and comparing also served as aids for the decision of annotation. The deviation of the pronunciation was then compared with the standard transcription and pronunciation from British lexicon,

4 Results

The results can be represented at word level and phoneme level.

4.1 Word Level

Each word can be represented in the following way (Figure 2), like that in Schaden [6]. The first level (word level) describes the orthographic form of the word and its standard transcription in SAMPA; the second level (phonetic variable level) illustrates the phonemes which can have potential deviations; the third level (variable value level) lists all the actual phone variants from different speakers. Only one speaker is listed here for one variable value, in practice we investigated all the speakers, and calculated the rate of different phone variants in our study.

For the lexical analysis, there are several levels of accent. We can take the word "Colorado" as an example in Table 2. Besides the information of accent level, the accent level rate is

Table 1 - English loan words in German SMS texts

Kelly-Clarkson-Album	Home	Colorado	PDA
Powerpoint-Folien	Date	Trainer	Playstation
Credit points	Sorry	To-do-list	Badmintonschläger
Broadway-Aufführung	House	Favorite	Explorer
Merry Christmas	Office	Bluetooth	Motor boot
Happy new year	Departed	iPod	Southampton
Harry Potter	Starclub	Superstar	Kingdom Come
Disneyland Resort Paris	Special	Blinddate	James Bond
Take me out to the ballgame	Software	Attachment	Coldplay-CD
All Inclusive	Controller	DiscoBoys	Big Brother
Travel and Work	Skateboard	Blood Diamond	See you
Prison Break	Clients	Action	Feed
Happy Birthday	Player	Township	Blood Diamond
San Francisco	E-Mail	The Gang	Extreme Activity
Counter Strike Source	Piercings	Dirty Dancing	Snowboardkurs
World of Warcraft	Everwood	Wrestler	Charity-Gala
Hardcore Snowboard	Bluetoothmodul	Airport	Rock and Roll
Little Miss Sunshine	Cowboystiefeln	Beachvolleyball	Sean Connery
Tower of London	Snowboardkurs	HomeZone	Chrysler Building
Salad dressing	Gilmore Girls	Outlet-Center	In My Mind
Streetball-Turnier	Supershow	Gecancelt	Football Book
Sneak Preview	Fantasy	Illustrated	Centralstation
Candlelightdinner	Sound	Daily	Dream Dance
InterCityExpress	Source-Code	Cornflakes	Entertainment
Newcomerband	Release	Details	Sightseeing
I'm sure about this	Compact	Cool	T-shirts
Oldtimerclub	Export		

also calculated, for example "Colorado" was read by 16 speakers, from Accent Level 0 to Level 4, the percentage of the pronunciation variants are 12.5%, 12.5%, 12.5%, 37.5%, 25% respectively. In Accent Level 4, besides the typical variations, all other untypical variations are also included, the phones in the brackets "()" are another realization of the phoneme, they can be more than one value.

Because we have more than hundred words, we can not illustrate each word in this paper in detail. To calculate the word accuracy rate only by judging whether all phonemes are the same as the standard is also meaningless, because the difference between accent level 1 and accent level 4 can not be regarded as the same, sometime it is impossible to ascertain the difference between accent level 1 and without accent only by perception of average listeners. Therefore the accuracy rate based only on comparing the phonemes at the word level is not explanatory, and will not be presented here.

4.2 Phoneme Level

It is perhaps reasonable to present the results at the phoneme level by calculating the realization of particular phoneme with different variations. There are too many phonemes to be presented here, we choose only two phonemes: one is vowel /@U/, the other is consonant /r/ at syllable

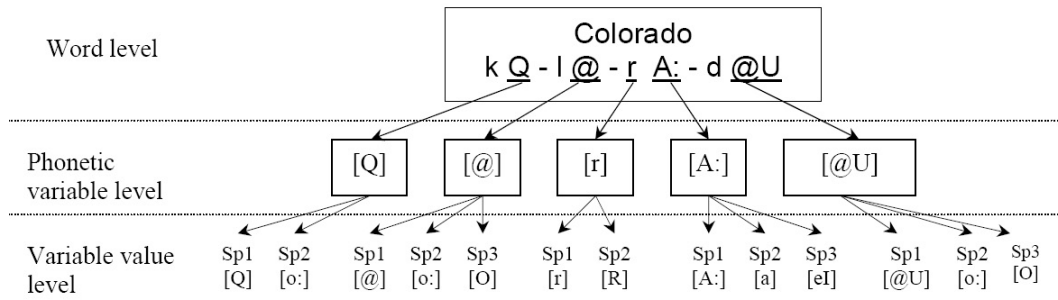


Figure 2 - Statistic scheme for phonetic variables

Table 2 - Accent level and accent rate of word "Colorado"

Accent Level	Transcription	Nr. of Speakers	Percentage
Without Accent	/ k Q - l @ - r A: - d @U /	2	12.5%
Level 1	/ k O - l @ - r a - d OU /	2	12.5%
Level 2	/ k o: - l o: - r a - d OU /	2	12.5%
Level 3	/ k o: - l o: - r a - d o: /	6	37.5%
Level 4	/ k o:(O) - l o:(O) - R a(eI) - d o:(O) /	4	25%

initial.

Actually it is still very difficult to give an average value of the realization, it depends on the word itself, if the word is already nativised, for example, /@U/ in "motor" would be generally pronounced as [o:], because it is generally accepted as a German word. Some speakers preferred [O] to [@U] in "controller", because they regarded it the same as "Kontroller". On the other hand, most of the speakers tried to pronounce [@U] in "home", but because of unprofessional English level, many failed in the realization of [@U] and replaced it with [OU] or [o:]. The realization of /@U/ is described in Figure 3. Statistics were based on 127 realizations of phoneme /@U/.

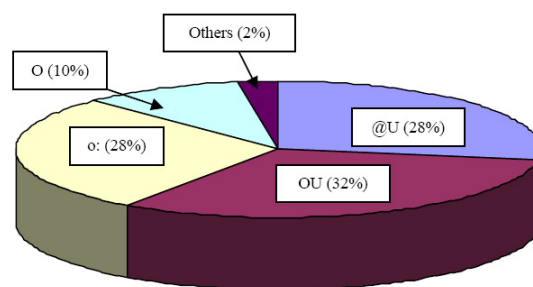


Figure 3 - Realization of /@U/ with different phonetic variables.

It is the same case with /r/ in Figure 4: Speakers preferred English [r] in not nativised words such as "prison" or "release"; Speakers were used to replacing [r] with German [R] in Germanized words such as "(Disneyland Resort) Paris" "(Salad) dressing". Statistics were based on 188 realizations of phoneme /r/.

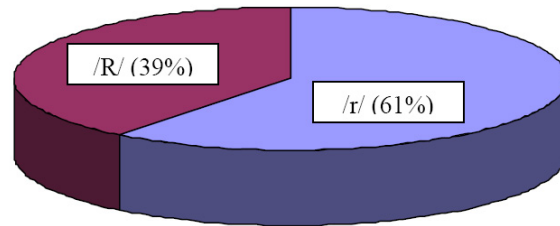


Figure 4 - Realization of /r/ with different phonetic variables.

5 Discussion

The results vary for the same phoneme, depending on the word it occurs in. Some influencing factors have to be mentioned here as being relevant: (1) age of the borrowing, (2) the way the loan words entered the German language, (3) general attitude to the word, (4) its origin etc. If the other conditions are the same, "partial substitution" loan words are likely to be pronounced according to German rules than "importation" words, for example, "Snowboard" in "Snowboardkurs" will be pronounced with more German accent than "Snowboard" in "Snowboard course".

Most of the English words contain proper names, such as place names or personal names, which deviate from the standard vocabulary, e.g. some speakers would pronounce "Colorado" in a wrong way as "/k o: - l o: - r eI - d o: /". The pronunciation of names are quite irregular in English, speakers have to rely their knowledge not only on the language but also on the information from TV, radio or communication with other people. Otherwise they will come across difficulties with the pronunciation of all kinds of names, especially the new coinages. The problems in proper name pronunciation is not only a problem of pronunciation with accent, but also concerns the knowledge of right pronunciation. On one hand, the deviation from the standard pronunciation of the names is quite difficult to predict, because their own pronunciations are irregular ; On the other hand, it is necessary to investigate the deviation, because they constitute a large part of the everyday communication, and should be faced by speech recognition and speech synthesis.

For speech synthesis purpose, it may be necessary to carry out a preference test to find out the most preferred variant, because it can be assumed that what the speakers produce themselves does not always resemble what they expect to hear from others [2]. The preferred variant can be of a "higher" level of pronunciation than the average speakers of German, which includes more English phonemes [2]. In the dictionary for synthesis, only one variation of pronunciation can be included, while in the dictionary of recognition, the possible variations of pronunciation should be taken into consideration in order to cover different accents of speakers. This paper only investigates the possible variation, and no preference tests have been conducted for the purpose of speech synthesis.

6 Conclusion

This concerned study investigated over 100 English words in SMS texts, 100 native German speakers were involved. Each word was uttered at least by 7 to 8 speakers (sometime twice in the texts), more than 1000 words were carefully annotated and statistically investigated. The words which have been studied can directly be used for the purpose of recognition of English words with German accent, the inferred rules at the word level and phoneme level can also be

used to predict the pronunciation variants of new English words. This investigation provides insight into the strategies of native German speakers in pronouncing different English words in SMS text domain. It further reveals which English phones are produced in which words, and to what extent, and which native German phones are used to substitute the English sounds in the communication of SMS text domain.

References

- [1] ABRESCH, J.: *Englisches in Gesprochenem Deutsch*. Dissertation, Universität Bonn, 2007.
- [2] ABRESCH, J. and S. BREUER: *Assessment of non-native phones in anglicisms by german listeners*. In *Proc. ICSLP 2004*, pp. 1281 – 1284, Korea, 2004.
- [3] BUSSE, U. and M. GÖRLACH: *German*. In GÖRLACH (ed.): *English in Europe*, pp. 13 – 36. Oxford: Oxford University Press, 2002.
- [4] DUCKWORTH, D.: *Zur terminologischen und systematischen Grundlage der Forschung auf dem Gebiet der englisch-deutsch Interferenz. Kritische Übersicht und neue Vorschlag*. In: KOLB, H. (Hrsg.): *Sprachliche Interferenz*, S. 36 – 56. Tübingen, 1977.
- [5] LANGNER, H. C.: *Die Schreibung englischer Entlehnungen im Deutschen*. Frankfurt am Main: Peter Lang, 1995.
- [6] SCHADEN, S.: *A database for the analysis of cross-lingual pronunciation variants of european city names*. In *Proc. LREC 2002*, pp. 1277 – 1283, Spain, 2002.
- [7] VELIU-AJDINI, S.: *Englisches in der deutschen Sprache der Werbung*. Diplomarbeit, Universität Wien, 2009.