



Diphthong modeling, phone mapping, and prosody transfer for speech synthesis of similar dialect pairs

Michael Pucher¹, Carina Lozo¹, Philip Vergeiner², Dominik Wallner²

¹Acoustics Research Institute, Austrian Academy of Sciences

²Department of German Studies, University of Salzburg, Austria

{michael.pucher, carina.lozo}@oeaw.ac.at

{philip.vergeiner, dominik.wallner}@sbg.ac.at

Abstract

Dialect synthesis is a challenging area of research and contrasts the synthesis of standard varieties not only as to the non standard nature of dialects but also in collecting proper data. In this paper we describe a diphthong interpolation and phone mapping based method that can be used to synthesize a new dialect with an existing dialect model of a similar dialect. The method only uses transcriptions of original dialect data, which are then mapped onto the phones in the model. We improve the basic mapping model further by transferring prosodic features such as original duration and F0. In addition to prosody transfer we want to investigate, if interpolation between two diphthong parts can substitute satisfactorily a missing phone in the target dialect. The methods are applied to two South-Bavarian dialects from Tyrol in Austria.

Index Terms: dialect, interpolation, prosody

1. Introduction

The collection of proper data for dialect synthesis is a time consuming task, due to the non-standard nature of dialects. As there is no standard orthography for dialects it is difficult to record dialect data as well as control a synthesizer through dialect input. Dialect speech synthesis is necessary for realizing different personas in spoken dialog systems since the perception of language variety (sociolect, dialect, accent) influences our evaluation of speaker's attributes (competence, intelligence, friendliness, etc.). These "dialect" personas are interesting for applications domains like regional information systems and computer games.

As we already showed in [1] dialect prosody has a positive influence on the authenticity of synthesized speech. In this paper, we extend the phone mapping method from [1] to two new target dialects and focus our modeling on diphthongs that are present in the target dialects but do not exist in the acoustic model. Therefore we investigate how interpolation between two diphthong parts can improve the synthesized output. With our methods we aim to synthesize a dialect of which no proper training corpus is available by using a trained acoustic Hidden-Markov-Model (HMM) of a similar dialect.

As our target dialects, we chose two dialects of Tyrol, the dialect of Schönwies (SWS), which is located in the political district of Landeck and Kelchsau (KLS), which belongs to the political district of Kitzbühel. The HMM for the source dialect was trained on data of the dialect of Innervillgraten (IVG) [2]. IVG presents the ideal source language, since it shares the super-ordinate dialect group (South Bavarian) with both SWS and KLS, and therefore also shares a similar phone set. Although KLS is a member of the South Central Bavarian and

SWS of the Bavarian Alemannic transition zone, due to the physical proximity to IVG we expect a beneficial commonality in their phone sets. We aim to use the similarity of these dialects to realize an authentic synthesized output for both the SWS and KLS dialects. Concerning the phone mapping, each SWS and KLS phone is separately mapped to a corresponding IVG phone, which defines a context-free mapping between the phone sets. The mapping is created manually but can in principle also be derived automatically with a vector distance approach, if phones are defined by feature vectors.

The phone sets of SWS and KLS were extracted via a short section of a longer recording. A total of only 20 utterances of authentic dialect speakers each from SWS and KLS were transcribed and were also used as test sentences for the synthesizer. For the evaluation we synthesized four different versions of the test sentences. The first sample is fully synthesized (*Syn*), the second is the basic synthesis with the interpolation (*Syn+diph*). The third sample includes the original SWS or KLS speaker's duration as well as the interpolation (*Syn+diph+dur*). The fourth version also includes the original speaker's segment duration and F0 (*Syn+diph+dur+F0*) along with the interpolation.

To evaluate the quality and the authenticity of the synthesized dialect, we distributed an online survey to the regions where SWS and KLS are located, reaching thereby only listeners who are familiar with Tyrolean dialects. With our approach, we achieve the synthesis of SWS and KLS exclusively by using the already created synthesis system for IVG and a phone mapping between SWS/KLS and IVG. By these means, we can prepare the ground for further studies and experiments involving poorly digitized dialects.

Previous work on the synthesis of language varieties investigated HMM-based synthesis of the modern Hanoi dialect [3], Austrian dialects [2], lexicon learning techniques with a smaller lexicon available for bootstrapping [4], data selection and front-end construction for different languages, varieties and speaking styles [5], and phoneme-to-phoneme conversion for converting pronunciations between American, British and South African English accents [6]. Voice model interpolation was applied in HMM-based synthesis for speaker interpolation [7], emotional speech synthesis [8], emphasis of foreign accent [9], interpolation of accented English within a reactive HMM-based synthesis system [10], and supervised interpolation between phonetically different dialects [11]. Prosody transfer was investigated for modeling of Swiss French accent [12].

2. Corpora

2.1. Dialects of Schönwies – Innervillgraten – Kelchsau

Both SWS and KLS belong to the South Bavarian dialect group, the same accounts for IVG. It must be stressed, however, that SWS as well as KLS are located within transition zones - SWS in the Bavarian Alemannic transition zone and KLS in the South Central Bavarian transition zone. In this sense they lack some South Bavarian dialect features. Examples of South Bavarian dialect features are the preservation of the postvocalic /l/ in South Bavarian in contrast to its vocalization in Central Bavarian (in KLS it is already vocalized) as well as the preservation of the distinction between standard and *Middle High German* (MHG) <ʰ> and <d> in contrast to its coincidence via lenition in Central Bavarian. Other characteristics of South Bavarian include the realization of falling diphthongs for MHG <ê, ô, œ> (e.g., [grœ̯s] for *groß*, “big”) and the occurrence of [kç] or [kh] for Germanic *k*. A typical dialect feature for the Tyrolean area shared with most Alemannic dialects is the palatalization of /s/ [13].

2.2. Data collection

As mentioned above, we decided that dialects from the same dialect group as IVG would benefit the experiment because of the expected overlaps in their phone sets. For this contribution data from three different corpora was used. The IVG, SWS, and KLS data was collected separately and within different time spans.

For the IVG corpus ten dialect speakers, gender balanced, were recruited. The recordings consisted of spontaneous speech, reading tasks, picture naming tasks, and translation tasks from Standard Austrian German (SAG) into the dialect. From these recordings, 660 phonetically balanced sentences were selected and a phone set was created for the IVG dialect. For the recording of the 660 phonetically balanced dialect sentences both the audio and the orthographic script, based on Standard German, of the samples to be collected were presented to the dialect speakers, who were then asked to repeat the individual samples in the IVG dialect. In addition, these speakers also read a corpus of SAG sentences. The speaker selection and recording process for IVG has been described in detail in [2]. A phone set for IVG was then created using this data, consisting of 82 phones. In this paper we use the data from the male IVG speaker LSC. 656 IVG dialect sentences were recorded from this speaker.

Sound samples were recorded at 44100 Hz, 16 bits/sample. The training process was also performed using these specifications. Cutting and selection was performed manually. Noise cancellation and volume normalization were applied to the recordings. Synthesized samples used in the evaluation were also volume normalized. A 5 ms frame shift was used for the extraction of 40-dimensional mel-cepstral features, fundamental frequency and 25-dimensional band-limited aperiodicity [14] measures. Speaker-dependent models were trained for the evaluations using the HSMM-based speech synthesis system published by the EMIME project [15].

The data used for SWS and KLS was collected by the FWF-project “German in Austria” [16]. Both for SWS and KLS four speakers were recorded, in each case two NORM/Fs (= non mobile, old rural males / females) and two young professionals (= speakers younger than 25, lacking upper-secondary education, born and working within the area). For the purpose of the present paper two informants, one NORM in SWS and one

male young professional in KLS, were selected. The recordings consisted of picture naming tasks and translation tasks from SAG into base dialect. The answers were elicited through a dialect survey by the same trained interviewer. In total 542 items were collected in each recording session which lasted up to 5 hours, sound samples were recorded at 44100 Hz, 16 bits/sample. From the translation tasks 20 sentences containing various diphthongs were chosen for the synthesis. These 20 sentences were then phonetically transcribed in IPA and SAMPA and manually segmented on the phone level with STx [17] and Praat [18].

3. Phone Mapping

The examined 20 sentences of both SWS and KLS allowed us to extract a small phone set for each dialect. On the basis of the phonetic transcription of 20 sentences, we determined 76 SWS and 57 KLS phones. The comparison between IVG and SWS, as well as IVG and KLS showed that 30% (n=22) of SWS’ and 30% (n=17) of KLS’ phone set are also included in the source dialect’s phone set. The target dialects’ phones without an exact match in the source dialect’s phone set were manually assigned to the most similar in the IVG phone set. Figure 1 illustrates the overlap between the source dialect IVG and both target dialects SWS and KLS. In the case of vowels, the mapped phone must be in a similar location of constriction and must not differ in more than two distinctive features (quality, dorsality, height, or roundedness) from the SWS/KLS vowel. For example, the close-mid front rounded vowel [ø] was mapped to the open-mid front vowel [œ], differing only in quality. As most of the mapped phones are diphthongs which differ in terms of tenseness from the IVG phones, we decided to map diphthongs in a similar way to vowels. For the IVG equivalents, we chose IVG diphthongs with the highest agreement in constriction location and the distinctive features mentioned above. By this approach, we could map 9 diphthongs, but some sounds still remained without a proper IVG counterpart. Since several diphthongs could not be mapped satisfactorily, we decided to interpolate between those parts. Tables 1 and 2 present the mapped phones, grey rows indicate interpolation. It has to be noted, that these tables only present the phones of the data which were subjectively evaluated by listeners, hence they do not include all mapped phones.

Table 1: *KLS mapping*. Table 2: *SWS mapping*. “*” signifies phone concatenation.

KLS	IVG	SWS	IVG
ĩ	ɪ	ɐ	ʏ
ɛ̃	ɛ̃	ç	ɐ̯
ɪ̃	ɪ̃	ɪ̃	ɪ̃
ũ	u * ɐ	ũ	u * ɐ
ö̃	ö * ə	ö̃	ö * ɐ
õ	o * ɐ	õ	õ
ç	ɐ̯	ɑ̃	ɑ * ʊ
ɐ	ʏ	ʊ	v
z	s		
v	v		

4. Synthesis

For synthesizing the four different versions of our prompts, the phonetic transcription of the KLS and SWS utterances were

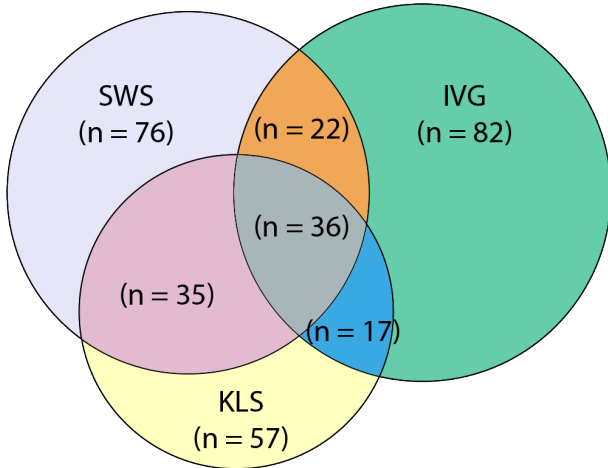


Figure 1: Overlap of SWS-IVG-KLS phone sets

transformed into phonetic transcriptions of IVG utterances using the mapping process described in Section 3. The mapped phones are shown in Table 1 and Table 2. This phonetic transcription which also contains syllable and word boundaries, was then transformed into full-context labels that can be used with the HMM-based synthesis system. The whole process emulates the function of a speech synthesis front-end that performs text normalization and Letter-to-Sound (LTS) conversion.

5. Prosody transfer

The process of prosody transfer is done in the same way as in our previous paper [1] and is only described shortly here. Using the full context labels and a speaker dependent voice of the IVG speaker LSC that we have developed, we synthesize the first set of prompts for the speaker KLS and SWS. The second set of prompts uses the phone duration information from the original KLS and SWS speakers.

As a third test set we additionally used the original speaker’s F0 values during synthesis, where we modified the SWS and KLS speaker’s original F0 curves $F0_{orig}$ as follows:

$$F0_{syn,dur,f0} = F0_{orig} + \left(\frac{\sum F0_{syn,dur}}{N} - \frac{\sum F0_{orig}}{N} \right). \quad (1)$$

N is the number of F0 values, e.g., frames in the utterance, i.e., this adaptation was done on the utterance level. $F0_{orig}$ is the original F0 trajectory. $F0_{syn,dur}$ is the synthesized F0 trajectory with original durations. Additionally, we used the voicing decision from the synthesized F0 curve.

The F0 algorithm was applied to only take over the F0 dynamics from the original speaker and to keep the F0 range within the range of the synthetic speaker, such that the comparison of synthesized samples is easier.

6. Diphthong modeling

In addition to the phone mapping and prosody transfer we also used diphthong modeling in this paper. As can be seen in Table 1 and 2 there are several diphthongs in the KLS and SWS variety that do not exist in the IVG variety that was used for the acoustic model. To be able to synthesize these diphthongs we applied two methods.

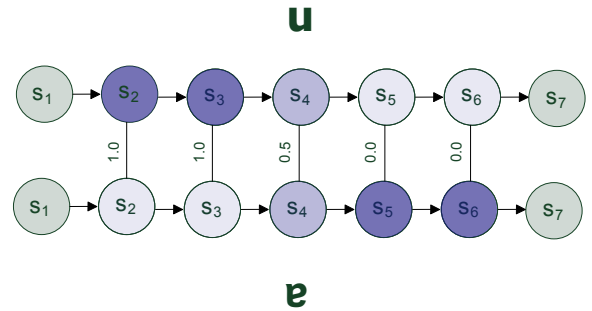


Figure 2: Diphthong interpolation between phone [u] and [v] to generate diphthong [uɐ]. Coloring indicates the weight of the node given by the interpolation.

6.1. Diphthong splitting

In the first method we split up the diphthong in KLS or SWS into two phones in IVG. The duration of the diphthong from was evenly distributed to the two mapped phones. A disadvantage of this method is that we synthesize a diphthong from two phones that have different right (for the first phone) and left (for the second phone) contexts. In this way we are not able to correctly model the transition from the first to the second phone that is happening in a diphthong.

6.2. Diphthong interpolation

To be able to better model the transition from the first to the second phone in the diphthong we implemented diphthong interpolation. The process of diphthong interpolation is shown in Figure 2 for two 5-state HMMs for the phones [u] and [v].

For the first two emitting states of the interpolated phone we only use the information of the first phone (interpolation ratio $\alpha = 1.0$), for the third state we use an interpolation of the first and second phone (interpolation ratio $\alpha = 0.5$), and for the last two states we use only the second phone (interpolation ratio $\alpha = 0.0$). In this way we do not use the two rightmost states of the first phone and the two leftmost states of the second phone, both being trained in inadequate contexts for the diphthong.

Mel-cepstral, F0, and aperiodicity parameters are interpolated linearly with

$$\mu = \alpha\mu_1 + (1 - \alpha)\mu_2 \quad (2)$$

$$\sigma^2 = \alpha^2\sigma_1^2 + (1 - \alpha)^2\sigma_2^2 \quad (3)$$

being the formulas for interpolated mean μ and variance σ^2 parameters of normally distributed random variables respectively. For the duration interpolation the mean values μ_1 and μ_2 are just added up to take into account the fact that we have no information on the diphthong duration and that diphthongs are normally longer than the monophthongs from which we construct them here.

The interpolation always happens within the context of an utterance. For the interpolation we first construct two phonetic input sequences where the first one contains the first phone of the diphthong and the second contains the second phone of the diphthong. These two sequences can be generated using the phone mapping from Section 3. Table 3 shows the two input sequences. The algorithm then looks for the differences in the input sequences and then applies the interpolation shown in Figure 2 to the 2 different models, the model for [u] and [v] in the [χ.st] context in this example.

Table 3: *Phonetic input/output sequences of interpolation.*

KLS orth.	Lauf so schnell du kannst.
KLS phon. input 1	laf sə ʃnœ̃ ðʊ χust
KLS phon. input 2	laf sə ʃnœ̃ ðʊ χɛst
KLS phon. output	laf sə ʃnœ̃ ðʊ χuɛst

7. Evaluation

For the subjective evaluation we used four different synthesis methods:

1. *Syn* - For this method only the phone mapping was used and synthesis was done from speaker-dependent HMMs. Diphthong splitting was used for diphthong modeling.
2. *Syn+diph* - This method used diphthong interpolation as described in Section 6.2.
3. *Syn+diph+dur* - This method also used prosody transfer with the duration from the original speech samples.
4. *Syn+diph+dur+F0* - This method also used the original F0 values.

We did not include methods with prosody transfer that do not use diphthong interpolation, since these methods were already evaluated for a different dialect previously [1].

The synthesized output for SWS and KLS was evaluated separately. For the SWS evaluation we had 7 (5 ♀, 2 ♂) listeners in total, 6 from the administrative district Landeck, 1 from the political district Innsbruck-Land in the province Tyrol, Austria. The KLS output was evaluated by a total of 21 listeners (9 ♀, 12 ♂), 17 from the political districts Kitzbühel, 2 from Kufstein, and 1 each from Landeck and Imst.

For the evaluation an online survey was set up with [19]. The evaluation consisted of a pairwise comparison of the different synthesis methods. For the subjective evaluation we chose 4 out of the 20 test sentences from the SWS and KLS corpus¹. These utterances were then synthesized in each of the four methods, hence making $4 * 4 = 24$ samples for KLS and SWS. Considering 6 possible combinations of the 4 methods, we presented the listeners with 24 pairwise comparisons, where they were asked to choose the better sample.

Figure 3 shows the result of the pair-wise comparison. According to a Wilcoxon rank sum test significant differences can be found between the methods *Syn* and *Syn+diph+dur+F0* for SWS and between *Syn* and *Syn+diph+dur+F0* as well as *Syn+diph* and *Syn+diph+dur+F0* for KLS. These results show that the basic synthesis method with phone mapping and diphthong splitting (*Syn*) achieves a reasonable modeling of diphthongs and duration, since there are no significant differences between *Syn*, *Syn+diph*, and *Syn+diph+dur*.

8. Discussion

Although this result might be disappointing from a purely speech synthesis point of view it is interesting from a dialectological viewpoint. Combined with our previous results on phone mapping and prosody transfer experiments for a different Austrian dialect (STY) [1] where we saw significant improvements by using the original duration, we can conclude that prosodic features have either different prominence for different

¹The samples used in the listening test can be found on <https://speech.kfs.oeaw.ac.at/dialektsynssw10/>

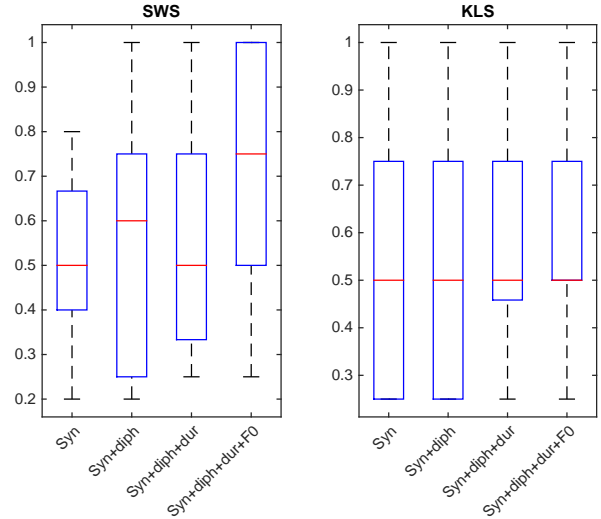


Figure 3: *Results of pairwise comparison.*

dialects or that the temporal organization patterns are an important feature of dialect proximity, since IVG is closer to KLS and SWS, than to STY. It also seems that the correct modeling of F0 patterns is of utmost importance for realistic dialect modeling.

Concerning diphthong interpolation we did not see an improvement with the dialects KLS and SWS, but it still may be a useful technique for dialect pairs that have a larger set of different diphthongs. Furthermore we see that the basic diphthong modeling method can achieve a reasonable accuracy of diphthongs, although diphthongs are thought as being a salient feature of KLS and SWS dialects.

9. Conclusion

In this paper we showed how diphthong interpolation, phone mapping and prosodic transfer can be used for the synthesis of a new dialect of which only symbolic phonetic input is available. In the pairwise comparison task we found a significant difference of some models to the model that uses original duration and F0 (*Syn+diph+dur+F0*). This shows us that the baseline model works reasonably well in terms of duration modeling and modeling of diphthongs for the dialects of KLS and SWS, although we also see dialect-specific results concerning duration modeling of other dialects from previous studies.

We also showed how speech synthesis can be used as a tool in dialectology to investigate questions concerning prosodic distance between dialects. Our experiments show that either KLS and SWS are close to IVG in terms of durational patterns, or these patterns are less important for KLS and SWS.

The developed methods allow us to build a full moderate quality synthesizer for a new dialect without any speech training data of that dialect.

10. Acknowledgments

This work was supported by the Austrian Science Fund (FWF) project DiÖ - Deutsch in Österreich (I2539-G23).

11. References

- [1] M. Pucher, C. Lozo, and S. Moosmüller, “Phone mapping and prosodic transfer in speech synthesis of similar dialect pairs,” in *28th Conference on Electronic Speech Signal Processing*, Saarbrücken, Germany, 2017, pp. 180–185.
- [2] M. Toman, M. Pucher, S. Moosmüller, and D. Schabus, “Unsupervised and phonologically controlled interpolation of Austrian German language varieties for speech synthesis,” *Speech Communication*, vol. 72, pp. 176–193, 2015.
- [3] T. T. T. Nguyen, C. d’Alessandro, A. Rilliard, and D. D. Tran, “HMM-based TTS for Hanoi Vietnamese: issues in design and evaluation,” in *Proceedings of the 14th Conference of the International Speech Communication Association (Interspeech 2013)*, Lyon, France, 2013, pp. 2311–2315.
- [4] N. Goel, S. Thomas, M. Agarwal, P. Akyazi, L. Burget, K. Feng, A. Ghoshal, O. Glembek, M. Karafiat, D. Povey, A. Rastrow, R. C. Rose, and P. Schwarz, “Approaches to automatic lexicon learning with limited training examples,” in *Proceedings of the 2010 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Dallas, Texas, USA, 2010, pp. 5094–5097. [Online]. Available: <http://dblp.uni-trier.de/db/conf/icassp/icassp2010.html#GoelTAABFGGKPRRS10>
- [5] O. Watts, A. Stan, R. Clark, Y. Mamiya, M. Giurgiu, J. Yamagishi, and S. King, “Unsupervised and lightly-supervised learning for rapid construction of TTS systems in multiple languages from ‘found’ data: evaluation and analysis,” in *Proceedings of the 8th ISCA Workshop on Speech Synthesis (SSW)*. ISCA, 2013, pp. 101–106.
- [6] L. Loots and T. Niesler, “Automatic conversion between pronunciations of different English accents,” *Speech Communication*, vol. 53, no. 1, pp. 75–84, Jan. 2011. [Online]. Available: <http://dx.doi.org/10.1016/j.specom.2010.07.006>
- [7] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura, “Speaker interpolation for HMM-based speech synthesis system,” *Journal of the Acoustical Science of Japan (E)*, vol. 21, no. 4, pp. 199–206, Jul. 2000.
- [8] M. Tachibana, J. Yamagishi, T. Masuko, and T. Kobayashi, “Speech synthesis with various emotional expressions and speaking styles by style interpolation and morphing,” *IEICE Transactions on Information and Systems*, vol. E88-D, no. 11, pp. 2484–2491, nov 2005.
- [9] M. L. G. Lecumberri, R. Barra-Chicote, R. P. Ramón, J. Yamagishi, and M. Cooke, “Generating segmental foreign accent,” in *Proceedings of the 15th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2014, pp. 1302–1306.
- [10] M. Astrinaki, J. Yamagishi, S. King, N. D’Alessandro, and T. Dutoit, “Reactive accent interpolation through an interactive map application,” in *Proceedings of the 14th Conference of the International Speech Communication Association (Interspeech 2013)*, F. Bimbot, C. Cerisara, C. Fougieron, G. Gravier, L. Lamel, F. Pellegrino, and P. Perrier, Eds. Lyon, France: ISCA, 2013, pp. 1877–1878. [Online]. Available: <http://dblp.uni-trier.de/db/conf/interspeech/interspeech2013.html#AstrinakiYKDD13>
- [11] M. Pucher, D. Schabus, Y. Yamagishi, F. Neubarth, and V. Strom, “Modeling and interpolation of Austrian German and Viennese dialect in HMM-based speech synthesis,” *Speech Communication*, vol. 52, pp. 164–179, 2010.
- [12] P.-E. Honnet, A. Lazaridis, J.-P. Goldman, and P. N. Garner, “Prosody in swiss french accents: Investigation using analysis by synthesis,” 2014. [Online]. Available: <http://infoscience.epfl.ch/record/198138>
- [13] P. Wiesinger, *The Central and Southern Bavarian Dialects in Bavaria and Austria*. Routledge, 1990, pp. 438–519.
- [14] H. Kawahara, I. Masuda-Katsuse, and A. Cheveigné, “Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based F0 extraction: possible role of a repetitive structure in sounds,” *Speech Communication*, vol. 27, pp. 187–207, 1999.
- [15] J. Yamagishi and O. Watts, “The CSTR/EMIME HTS system for Blizzard challenge 2010,” in *Proceedings of the Blizzard Challenge Workshop*, Kyoto, Japan, 2010, pp. 1–6.
- [16] G. Budin, S. Elspaß, A. Lenz, S. M. Newerkla, and A. Ziegler, “The research project (SFB) ‘German in Austria’. Variation - Contact - Perception,” *Zeitschrift für germanistische Linguistik. Deutsche Sprache in Gegenwart und Geschichte. Dimensions of Linguistic Space: Variation – Multilingualism – Conceptualisations*, vol. 46/2, pp. 7–35, 2019.
- [17] A. Noll, White, J., Gottschall, C., Becker, T., and Balazs, P., “Sound Tools Extended (STx) version 4.1.0,” <http://www.kfs.oeaw.ac.at/STx/>, 2014.
- [18] P. Boersma and D. Weenink, “Praat: doing phonetics by computer. version 6.0.37,” <http://www.fon.hum.uva.nl/praat/>, 2018.
- [19] D. Leiner, “SoSci Survey. Version 3.1.06–i,” <http://www.sosicurvey.de/>, 2019.