# SwissDial: Parallel Multidialectal Corpus of Spoken Swiss German

*Pelin Dogan-Schönberger,   Julian Mäder,   Thomas Hofmann*

ETH Zürich, Department of Computer Science, Switzerland

{pelin.dogan, maederju, thomas.hofmann}@inf.ethz.ch

## Abstract

Swiss German is a dialect continuum whose natively acquired dialects significantly differ from the formal variety of the language. These dialects are mostly used for verbal communication and do not have standard orthography. This has led to a lack of annotated datasets, rendering the use of many NLP methods infeasible. In this paper, we introduce the first annotated parallel corpus of spoken Swiss German across 8 major dialects, plus a Standard German reference. Our goal has been to create and to make available a basic dataset for employing data-driven NLP applications in Swiss German. We present our data collection procedure in detail and validate the quality of our corpus by conducting experiments with the recent neural models for speech synthesis.

**Index Terms**: dataset, speech synthesis, low-resource, NLP, Swiss German dialects

## 1. Introduction

In recent years, research in natural language processing has led to significant advancements on tasks such as sentiment analysis, question-answering, text- and speech-based machine translation, automatic speech recognition, or speech synthesis. Intelligent language systems such as voice assistants, chat bots, or AR/VR devices start to permeate many aspects of our daily lives. These novel forms of human-computer interaction require suitable content and crucially depend on the availability of powerful models along with adequate data to train them. This creates inequalities across different languages and particularly poses challenges for low-resource languages and dialects, where data of sufficient quantity and quality is often scarce.

Among the low-resources languages, Swiss German dialects are characterized by being derived from Standard German (High German, DE) with considerable differences in phonetics, vocabulary, morphology and syntax, even lacking clear boundaries and definitions. Unlike other languages, where dialect usage decreases in favor of standard variants in most social domains, Swiss German dialects are widely used in everyday life (cf. [1]). With the expansion of personalized digital media and social platforms, dialect use is growing even more. Despite this increasing demand, the obstacles posed by, for instance, lack of standard orthography, have made data collection of written text (e.g. transcripts) challenging. These difficulties in collecting consistent and clean data have led to a situation, where Switzerland is lagging behind and struggles to take full advantage of basic NLP tools.

This paper presents an annotated parallel corpus of spoken Swiss German dialects, *SwissDial*, for 8 different regions: Aargau (AG), Bern (BE), Basel (BS), Graubünden (GR), Luzern (LU), St. Gallen (SG), Wallis (VS) and Zürich (ZH). Our corpus consists of web-crawled sentences in High German, manual translations into Swiss German dialects and their audio recordings. Different genres and sources are represented: news sto-



```
{  "id": 57,
   "topic": "special",
   "code-switching": false,
   "DE" : "Dann lacht er sie aus und springt fort.",
   "AG" : "Dene heter sie usglachet ond esch weggloffe.",
   "BE" : "De lachtersi us und schpringt furt.",
   "BS" : "Denn lacht är si us und springt wäg.",
   "LU" : "Denn lachter sie us ond rennt devo.",
   "SG" : "Denn lacht er sie us und springt los.",
   "GR" : "Denn lacht er si uus und springt drvoo.",
   "VS" : "De lacht är schi üs und springt ewäg.",
   "ZH" : "Dänn lacht er si us und rännt furt."  }
```

Figure 1: *A sample sentence from SwissDial. For each sample sentence in High German, we provide the manual translations and the speech samples in 8 Swiss dialects. (The sentence means "Then he laughs at her and jumps away." in English)*

ries, Wikipedia articles, weather reports, short stories. The audio recordings are performed by a single speaker for each dialect. Figure 1 presents a sample from SwissDial. Our primary motivation and focus is to provide a basic dataset that would allow applicability of recent data-driven NLP models, especially for speech synthesis, and step up the ongoing research in Swiss German. To our knowledge, this is the first publicly available parallel corpus of spoken Swiss German. Our main contributions are as follows:

- We present *SwissDial*, a parallel corpus of text and audio in 8 Swiss German dialects on designated topics. The data collection process is well-documented, allowing for transparency and augmentability. The corpus is available at https://projects.mtc.ethz.ch/swiss-voice-data-collection.

- We validate the suitability of our corpus by running experiments with the recent neural models for (i) single speaker, (ii) multi-speaker-multi-dialect data, and (iii) code-switching (CS) speech synthesis. The samples from these models are available at https://projects.mtc.ethz.ch/projects/swiss-voice/swissdial.

## 2. Related Work

For various high-resource languages, there are parallel text corpora [2, 3, 4] and monolingual text corpora such as [5], Google News, Common Crawl that are used for the tasks of language modeling and machine translation. [6, 7, 8] present paired audio and text data which are widely used for automatic speech recognition and speech synthesis with various transformations. For these languages, [9] presents a parallel corpus for multilingual speech translation.

However, availability of data is limited for languages containing strong dialectal variations such as Arabic, Chinese, Swiss German which come forward with their special condition where the formal variety of the language differs significantly from varieties that are acquired natively and mostly used for verbal communication. There are various datasets concerning Swiss German that allow various NLP tasks such as dialect identification [10, 11], dialect machine translation [12, 13], morphology generation [14], automatic speech recognition [15]. Archimob [16, 17] presents a general purpose corpus of spoken Swiss German based on oral history interviews with various people from different dialects, which provides pairs of audio and text with annotated normalization and part-of-speech tagging. Another resource [18], Swiss SMS corpus, provides text corpus of SMS messages. NOAH's corpus [19] presents a collection of text in various genres in Swiss dialects with manually annotated part-of-speech tags. [12] collects a various parallel written sources (High German and Swiss German) to perform machine translation by normalizing Swiss German. Recently, [20] introduced a dictionary containing forms of common words in various Swiss German dialects normalized into High German.

Most of these existing resources in Swiss German do not provide sufficient data to train end-to-end models for neural machine translation and speech synthesis. They rather provide word-level correspondences which are suboptimal for sentence level NLP due to the lack of full context and syntax. Moreover, existing datasets of spoken language do not fulfill the requirements of today's neural speech synthesizers, which ideally need clean and single speaker audio. SwissDial is thus, to our knowledge, the first available parallel multi-dialectal corpus with transcribed text and clean audio in Swiss German. The text part of our corpus is carefully designed to cover a large set of topics and lexicons providing manual and parallel translations of sentences in dialects. The audio part, paired with the corresponding text, provides clean and single speaker audio for each dialect at sentence level. Moreover, the corpus contains labeled code-mixing samples which would be valuable for code-mixing speech synthesis research. We believe that the parallel structure of our corpus can also be used to improve data efficiency in NLP tasks via transfer learning.

## 3. SwissDial Corpus Content and Creation

In the following, we present the construction and content of our corpus and its modalities. Acquisition of SwissDial for applications in end-to-end neural models has five main steps: obtaining raw data, selection of dialects and annotators, translation, recording, and post-processing.

### 3.1. Raw data

Swiss dialects are spoken colloquially, but only rarely used in written form. Moreover, the available written forms collected from unspecific sources generally contain inconsisten-

cies within the same dialect, since there is no official standard orthography. This poses a challenge for annotators, possibly leading to disagreements with regard to the orthography of the provided text. To increase consistency and reduce ambiguity of the process, we started by collecting High German sentences to be translated by the annotators themselves.

We collected random sentences from news articles to cover a wide range of topics for good generalization. These sentences were labeled with the topic metadata of the article that they belong to. As mentioned earlier, Swiss dialects differ in vocabulary from High German and from each other. We thus crawled the internet for lists of lexical items that show variability among the dialects. We then extracted sentences from Wikipedia articles that contain each word (labeled as *special*). In this way, we aimed to cover vocabulary differences of dialects, while having a wide range of topics for a good generalization in various NLP tasks. Table 1 presents the topic distribution of the collected sentences in the corpus. It is important to note that significant majority of the code-switching (CS) content was the result of the random selection procedure, rather manual addition. There are two main reasons for this relatively large amount of CS content: (i) there is a high linguistic contact between English and German which results in considerable amount of word *borrowing*, (ii) our main sentence source is news articles that contain global topics like politics, science, technology etc, and therefore naturally contain words from English.

### 3.2. Selection of dialects

Swiss dialects are not clearly separated by geological borders, the variation in the regional language is rather continuous. One can separate each dialect into numerous sub-dialects even down to resolution of villages. Therefore, we had to make discretization where we aimed to represent the differentiation of the dialects as much as possible while addressing the most populated areas. Our selection motivation was to cover a range of dialects while complying with time constraints and availability of annotators. As a result, our corpus contains data from the following regions: Aargau, Bern, Basel, Graubünden, Luzern, St. Gallen, Wallis and Zürich. To represent these areas, we carried out an auditioning where we evaluated the annotator candidates based on their performance on a sample translation and pronunciation task, as well as their vocal quality features inspired by [21]. We hired one non-professional annotator for each dialect that passed all the tests to perform the translation and audio recording steps.

### 3.3. Text translation

We asked annotators from each dialect to translate the provided High German sentences, described in Section 3.1, to their native dialect. As mentioned earlier, Swiss dialects have different morphology together with different vocabulary and word choice. In some cases, a High German sentence or expression may not be easily translatable into dialect without sounding artificial. Therefore annotators had the freedom to skip such sentences. All the numbers, abbreviations and time expressions were translated as the annotator would read them out loud in the next step. In other words, the translation step aimed to have a one-to-one mapping between the text and the audio to-be-recorded. The foreign words, e.g. code-switching words, in the High German sentences were directly transferred as in the original language, rather than performing translation or phone mapping.

| Topic | Source | Parallel | AG | BE | BS | GR | LU | SG | VS | ZH | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Animals/farming | News | 130 | 137 | 135 | 137 | 137 | 135 | 137 | 137 | 134 | 1'089 |
| Culture | News | 81 | 85 | 87 | 87 | 86 | 86 | 87 | 87 | 86 | 691 |
| Earth/Space | Wikipedia | 17 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 25 | 151 |
| Economics | News | 60 | 63 | 63 | 62 | 63 | 63 | 63 | 63 | 234 | 674 |
| Int. politics | News | 140 | 155 | 154 | 153 | 155 | 152 | 155 | 155 | 151 | 1'230 |
| Medicine | News | 104 | 111 | 110 | 109 | 112 | 110 | 112 | 112 | 140 | 916 |
| Meteorology | News | 192 | 215 | 201 | 214 | 215 | 202 | 215 | 215 | 211 | 1'688 |
| Random | Wikipedia | 132 | 149 | 142 | 142 | 148 | 148 | 149 | 147 | 160 | 1'185 |
| Special | Wikipedia | 1030 | 1'120 | 1'105 | 1'106 | 1'120 | 1'112 | 1'122 | 1'123 | 1'730 | 9'538 |
| Science | News | 189 | 197 | 192 | 195 | 197 | 194 | 197 | 197 | 243 | 1'612 |
| Sports | News | 88 | 103 | 101 | 101 | 103 | 97 | 102 | 103 | 163 | 873 |
| Story | Internet | 49 | 50 | 50 | 50 | 49 | 50 | 50 | 50 | 81 | 430 |
| Swiss politics | News | 228 | 248 | 245 | 245 | 249 | 241 | 249 | 249 | 240 | 1'966 |
| Swiss regional | News | 88 | 97 | 97 | 94 | 97 | 97 | 96 | 97 | 467 | 1'142 |
| with Code-Switching | News | 232 | 247 | 247 | 247 | 249 | 248 | 250 | 248 | 345 | 2'081 |
| All topics | - | 2528 | 2'748 | 2'700 | 2'713 | 2'749 | 2'715 | 2'752 | 2'753 | 4'065 | 23'195 |

Table 1: *The distribution of sentences among dialects.* Special *indicates the sentences that are manually prepared to catch different vocabulary among the dialects.*

| Dialect | Total Audio (h) | Average utterance (s) | Average #word per utterance |
|---|---|---|---|
| AG | 2.78 | 3.64 | 11.8 |
| BE | 3.41 | 4.55 | 11.9 |
| BS | 3.15 | 4.18 | 12.7 |
| LU | 2.54 | 3.36 | 12.3 |
| GR | 2.93 | 3.83 | 12.5 |
| SG | 3.71 | 4.85 | 12.1 |
| VS | 3.32 | 4.34 | 11.9 |
| ZH | 4.55 | 4.03 | 12.2 |

Table 2: *The statistics of the collected audio data.*

### 3.4. Recording

As the last step of the data collection, the annotators were asked to read out loud their translations in a quiet room, in front of a high-quality recording set-up to obtain clean recordings. We used a Neumann TLM 102 microphone with pop filter and recorded in mono at 48 kHz. The speakers were instructed to speak at a distance of about 15 cm from the microphone. We prepared a user interface that shows the annotator the translated sentences each at a time and allows recording control through keyboard. In this way, the annotators were able to prepare and control their voice and breathing for each sentence, and rest when they needed. The interface also allowed annotators to modify their translations when they notice typos or incorrect wording/expressions before recording each sentence. This can be considered as the last quality control step for the correctness of the translations. The recordings were carried out in multiple sessions, around the same time of the day, that took 1-2 hours each. This scheduling is important to have a consistent voice across all the audio samples without exhausting the annotators and their vocal chords. Table 2 presents the statistics of the collected text data.

### 3.5. Post-processing

After translations and recording, we performed small post-processing on the collected data to make it more convenient

for various NLP tasks. For the text data in High German, we post-processed each sentence in the corpus by spelling numeric values into words. For the dialect text, we post-processed each sentence transforming the spelling of the numbers into numeric characters. In this way, we provide two versions for each sentence in all dialects: text with numeric characters and all characters. Providing both versions is helpful to tokenize numbers in various NLP tasks.

For the audio recordings, we performed a cleaning procedure where the aim was to cut off the segments at the beginning and end of the utterances that contain various noises like key pressing, breathing, mouth noises, etc. Some of the utterances contained noise during the sentence reading, rather than beginning/end, which would not be straightforward to clean. We repeated the recording step for these samples whenever it was possible, and eliminated the rest to present clean data which is significant for the speech synthesis task. Then, we down-sampled all recordings to 22050 Hz to be used in our speech synthesis experiments.

## 4. Speech Synthesis Experiments with SwissDial

With the increase in monolingual corpora, the end-to-end architectures [22, 23] have shown high quality speech synthesis results. Using these methods, we performed various experiments to synthesize speech in Swiss dialects using audio-text pairs from Swiss Dial. All the experiments below are carried out with character level input representations.

### 4.1. Single speaker model

We implemented the well-established neural speech synthesis model [22], an encoder-decoder architecture with attention, and conducted experiments on each dialect separately. To increase the data efficiency for our low-resource set-up, we firstly trained on the German part of CSS10 [8] (16 hours) for a good initialization. The main motivation for pretraining on High German is to transfer the textual and acoustic representations and learned alignment from High German data to Swiss dialects in the tuning stage, since Swiss dialects are derived from High German. Then, we fine-tuned the pretrained model on each dialect sep-
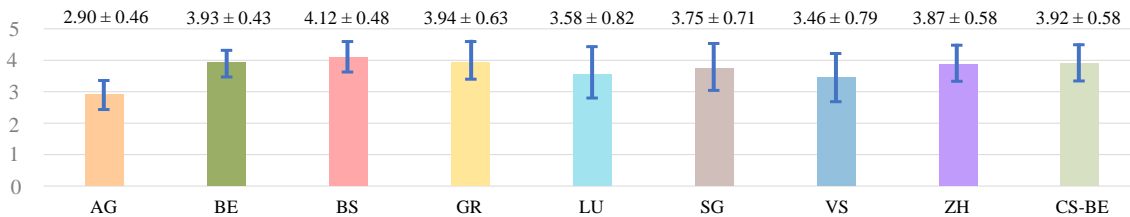
Figure 2: *MOS results on the single speaker model and code-switching model with standard deviation.*

arately to obtain a single speaker speech synthesis model. For each dialect, we held out 10% of the samples as a validation set.

We conducted a 5-scale mean opinion score (MOS) test, with 0.5 increments, for speech naturalness and quality. Fifteen native Swiss speakers from different regions were invited to participate. Randomly chosen ten utterances from a test set are presented to the participants for each dialect. Figure 2 shows the mean and deviation of MOS ratings of the participants. One has to keep in mind, that we hired non-professional speakers which naturally have, in comparison to professional speakers, much stronger deviations in quality and consistency of their pronunciation, intonation and voice control. Therefore, we suspect the differences in MOS rating between the dialects to be a result of the varying vocal quality of our non-professional speakers.

### 4.2. Multi-speaker-multi-dialect model

Next, we aimed to encourage sharing of model capacity across different dialects by training a single model on all dialects simultaneously. We implemented the multi-speaker model in [24], which is an extension of the previous model, and extended it with learnt dialect embeddings to support multi-dialect training. The speaker and the dialect embeddings are learnt in a similar way following [25]. We trained the speaker and dialect verification models by collecting samples from additional speakers from radio shows to enhance the capability of representations, since our dataset would only provide 8 different speakers. We held out 5% of the samples as a validation set. Different from [24], the learnt speaker and dialects embeddings are concatenated at the encoder output and the decoder input inspired by [26]. An internal subjective comparison between the multi-speaker-multi-dialect model results and the single speaker model results revealed a very similar naturalness and quality for the two models.

### 4.3. Code-Switching.

Lastly, we implemented the model *SE-DEC* [26], which is again an extension on [22], to explore the performance of our corpus in code-switching set-up with English words mixed in single Swiss dialect, BE, for showcasing. Firstly, we pretrained the model with monolingual recordings in English [7] and High German [8]. Then we fine-tuned the model using monolingual English data and the BE part of our dataset which introduces 247 code-switching samples. The MOS result for this experiment is shown in Figure 2 with the CS-BE bar. The participants were asked to rate the samples according to (i) Swiss German pronunciation, (ii) English pronunciation, (iii) naturalness and quality of the audio. As one can see in Figure 2, the resulting MOS ratings for BE with and without code-switching are almost identical, which is a very encouraging result for the use of code-switching extensions for Swiss German speech synthesis.

## 5. Conclusion

In this paper, we present SwissDial, a parallel multidialectal corpus of text and audio in 8 Swiss German dialects, by describing the collection procedure in detail. Furthermore, we validate the quality of SwissDial by presenting the results of experiments with the recent neural models for speech synthesis in single-speaker-single-dialect, multi-speaker-multi-dialect, and code-switching set-ups. Last but not least, we believe that applications on SwissDial are not limited only to speech synthesis but it can also empower further NLP research and application in Swiss German in the fields of machine translation and dialect identification.

## 6. References

[1] G. Liidi, "The Swiss model of plurilingual communication," *Receptive Multilingualism. Amsterdam: John Benjamins*, pp. 159–178, 2007.

[2] B. Jawaid, A. Kamran, and O. Bojar, "A Tagged Corpus and a Tagger for Urdu." in *LREC*, 2014, pp. 2938–2943.

[3] P. Koehn, "Europarl: A parallel corpus for statistical machine translation," in *MT summit*, vol. 5. Citeseer, 2005, pp. 79–86.

[4] C. Buck, K. Heafield, and B. V. Ooyen, "N-gram Counts and Language Models from the Common Crawl." in *LREC*, vol. 2. Citeseer, 2014, p. 4.

[5] M. Kupietz, C. Belica, H. Keibel, and A. Witt, "The German Reference Corpus DeReKo: A primordial sample for linguistic research," in *LREC*, 2010.

[6] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an ASR corpus based on public domain audio books," in *IEEE International Conference on Acoustics, Speech and Signal processing*, 2015, pp. 5206–5210.

[7] K. I. and L. J., "The LJ Speech Dataset," https://keithito.com/LJ-Speech-Dataset/, 2017.

[8] K. Park and T. Mulc, "CSS10: A collection of single speaker speech datasets for 10 languages," *arXiv preprint arXiv:1903.11269*, 2019.

[9] M. A. Di Gangi, R. Cattoni, L. Bentivogli, M. Negri, and M. Turchi, "MuST-C: A Multilingual Speech Translation Corpus," in *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 2012–2017.

[10] M. Ali, "Character level convolutional neural network for German dialect identification," in *VarDial*, 2018, pp. 172–177.

[11] T. Jauhiainen, H. Jauhiainen, and K. Lindén, "HeLI-based experiments in Swiss German dialect identification," in *VarDial*, 2018.

[12] P. E. Honnet, A. Popescu-Belis, C. Musat, and M. Baeriswyl, "Machine translation of low-resource spoken dialects: Strategies for normalizing swiss german," in *LREC 2018*, 2018.

[13] Y. Scherrer and N. Ljubešić, "Automatic normalisation of the Swiss German ArchiMob corpus using character-level machine translation," in *Proceedings of the 13th Conference on Natural Language Processing*, 2016.

[14] Y. Scherrer, "Morphology generation for Swiss German dialects," in *International Workshop on Systems and Frameworks for Computational Morphology*. Springer, 2011, pp. 130–140.

[15] P. N. Garner, D. Imseng, and T. Meyer, "Automatic Speech Recognition and Translation of a Swiss German Dialect: Walliserdeutsch," in *Proceedings INTERSPEECH 2014 – Annual Conference of the International Speech Communication Association*, 2014.

[16] T. Samardzic, Y. Scherrer, and E. Glaser, "Archimob-a corpus of spoken swiss german," in *LREC*. European Language Resources Association, 2016.

[17] Y. Scherrer, T. Samardžić, and E. Glaser, "Digitising Swiss German: how to process and study a polycentric spoken language," *Language Resources and Evaluation*, vol. 53, no. 4, pp. 735–769, 2019.

[18] E. Stark, B. Ruef, and S. Ueberwasser, "Swiss SMS corpus," https://sms.linguistik.uzh.ch, 2009-2015.

[19] N. Hollenstein and N. Aepli, "Compilation of a Swiss German dialect corpus and its application to PoS tagging," in *Proceedings of the First Workshop on Applying NLP Tools to Similar Languages, Varieties and Dialects*, 2014, pp. 85–94.

[20] L. Schmidt, L. Linder, S. Djambazovska, A. Lazaridis, T. Samardžić, and C. Musat, "A Swiss German dictionary: Variation in speech and writing," *arXiv preprint arXiv:2004.00139*, 2020.

[21] G. B. Kempster, B. R. Gerratt, K. V. Abbott, J. Barkmeier-Kraemer, and R. E. Hillman, "Consensus Auditory-Perceptual Evaluation of Voice: Development of a Standardized Clinical Protocol," *American Journal of Speech-Language Pathology*, vol. 18, pp. 124–132, 2009.

[22] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. Skerrv-Ryan *et al.*, "Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 4779–4783.

[23] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals *et al.*, "Wavenet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016.

[24] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, J. Shen, F. Ren, Z. Chen, P. Nguyen, R. Pang, I. L. Moreno *et al.*, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 4485–4495.

[25] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, "Generalized end-to-end loss for speaker verification," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 4879–4883.

[26] L. Xue, W. Song, G. Xu, L. Xie, and Z. Wu, "Building a Mixed-Lingual Neural TTS System with Only Monolingual Data," in *Proceedings INTERSPEECH 2019 – Annual Conference of the International Speech Communication Association*, 2019, pp. 2060–2064.