# Universal Automatic Phonetic Transcription into the International Phonetic Alphabet

*Chihiro Taguchi[1], Yusuke Sakai[2], Parisa Haghani[3], David Chiang[1]*

[1]University of Notre Dame, United States
[2]Nara Institute of Science and Technology, Japan
[3]Google, United States

ctaguchi@nd.edu, sakai.yusuke.sr9@is.naist.jp, parisah@google.com, dchiang@nd.edu

## Abstract

This paper presents a state-of-the-art model for transcribing speech in any language into the International Phonetic Alphabet (IPA). Transcription of spoken languages into IPA is an essential yet time-consuming process in language documentation, and even partially automating this process has the potential to drastically speed up the documentation of endangered languages. Like the previous best speech-to-IPA model (Wav2Vec2Phoneme), our model is based on wav2vec 2.0 and is fine-tuned to predict IPA from audio input. We use training data from seven languages from CommonVoice 11.0, transcribed into IPA semi-automatically. Although this training dataset is much smaller than Wav2Vec2Phoneme's, its higher quality lets our model achieve comparable or better results. Furthermore, we show that the quality of our universal speech-to-IPA models is close to that of human annotators.

**Index Terms**: speech recognition, phonetics, natural language processing, language documentation

## 1. Introduction

It is estimated that there are approximately 7,000 languages in the world, at least half of which will be extinct by the next century [1]. To record and revitalize these endangered languages, there have been many efforts by field linguists to document them and provide grammatical sketches and dictionaries. However, manual transcription of the recorded audio consumes a large amount of time in the documentation process [2]. Given this situation, the technology of automatic speech recognition (ASR) has the potential to accelerate language documentation.

This paper presents a model that converts speech to the International Phonetic Alphabet (IPA), which is the standard used in phonemic and phonetic transcription of understudied languages. We introduce a simple yet high-performance speech-to-IPA model that outperforms the current state-of-the-art, Wav2Vec2Phoneme [3]. Like Wav2Vec2Phoneme, our model is based on wav2vec 2.0 [4] and is fine-tuned to predict IPA from audio input. We use training data from seven languages in the CommonVoice[1] dataset, transcribed into IPA semi-automatically; that is, the quality of IPA transcription tools is verified to be reliable enough. The experimental results show that our model with 1,000 training samples per language (about 9 hours in total) outperforms existing speech-to-IPA models trained on more samples with more languages. The contributions of this paper are summarized as follows:

- We developed a new state-of-the-art speech-to-IPA model with a small training dataset;

- The results suggest that the quality of phonetic IPA labels and the linguistic diversity in the training dataset are significant factors of the performance in this task;

- The models, datasets, and G2P tools developed for this study will be publicly available.[2]

## 2. Related Work

Multilingual ASR has been advancing towards the ultimate goal of universal transcription for all spoken languages [5, 6, 7, 8]. However, mainstream objectives of ASR models have been to transcribe speech into graphemes, which are dependent on each language and its orthography; in other words, they are largely language-dependent rather than language-universal. Since endangered languages tend not to have any established orthography, they are often marginalized from computational processing methods that are based on orthography.

Some efforts have been made to train multilingual ASR models to transcribe low-resource languages into phones or phonemes written in IPA [9, 10, 11, 12]. In particular, Michaud [13] has carried out a successful case study of applying ASR to the documentation of the Na language. In recent years, more work on multilingual speech-to-phone ASR has been investigated using neural network models. Allosaurus [14] is a speech-to-IPA tool that is aimed to transcribe more than 2,000 languages. Gao et al. [15] train a wav2vec-based model on about 1,500 hours of audio from 15 languages; the text transcriptions are converted to IPA using a grapheme-to-phoneme (G2P) system. Wav2Vec2Phoneme [3] is a wav2vec2-based model that has a similar approach to ours. It fine-tunes wav2vec2-large-xlsr-53 [5] on approximately 2,000 hours of audio from 26 languages from CommonVoice. A significant difference between our approach and these models is the preprocessing phase to generate labels in IPA. Their IPA transcription labels are fully automatically generated using rule-based G2P tools, which do not necessarily guarantee the quality of its phonetic or phonemic transcription, while we only used G2P tools that are verified to be reliable enough. The details of the preprocessing step are described in the next section.

## 3. Method

### 3.1. Data

Datasets containing audio with close transcriptions into IPA are extremely scarce. Linguistic resources with speech–IPA pairs such as Pangloss [16] often contain IPA transcriptions that are not phonetic but phonemic or use language-specific conventions. To overcome this difficulty, most existing speech-to-IPA

---

Table 1: *A description of the datasets used for the training and validation sets in the study. "# Train" in the middle column refers to the maximum size of the training set available in CommonVoice 11.0. Seven languages are picked from CommonVoice for training the model.*

| CommonVoice | # Train | G2P tool |
|---|---|---|
| Japanese (ja) | 6,150 | Manual rules |
| Polish (pl) | 15,419 | Epitran |
| Maltese (mt) | 1,780 | Manual rules |
| Hungarian (hu) | 7,044 | Manual rules |
| Finnish (fi) | 2,044 | Manual rules |
| Greek (el) | 1,722 | Manual rules |
| Tamil (ta) | 16,189 | Epitran |
| Total | 50,348 | |

Table 2: *A description of the test set for testing the model on unseen languages. We picked four typologically diverse low-resource languages that are not part of the training dataset. A trained human annotator transcribed randomly chosen audio samples up to 100 samples in total.*

| CommonVoice | # Samples | G2P tool |
|---|---|---|
| Luganda (lg) | 22 | Manual annotation |
| Upper Sorbian (hsb) | 24 | Manual annotation |
| Hakha Chin (cnh) | 25 | Manual annotation |
| Tatar (tt) | 29 | Manual annotation |

models [3, 14, 15] use automatic G2P tools to transliterate the text of the datasets into IPA. However, because the output of these G2P tools is phonemic, and their quality is not guaranteed, it is possible that the generated IPA transcriptions do not represent the actual pronunciation accurately.

In this study, to alleviate this problem, we semi-automated the IPA transcription in two ways. First, we manually checked the quality of the G2P tools and only used those that were reliable enough. Second, to leverage the utility of the G2P tools as much as possible, we particularly chose languages that have a consistent orthography-to-pronunciation mapping. Based on these policies, we picked seven languages (Japanese, Polish, Maltese, Hungarian, Finnish, Greek, and Tamil) from CommonVoice 11.0.[3] For G2P, we used a combination of Epitran[4] [17] and tools implemented by us. Non-phonetic characters including punctuation symbols and spaces were removed at this point. Also, tones and other suprasegmental elements were not included in the transcription, because describing all the pitch information for each syllable of any languages would have yielded too much unnecessary information for further linguistic analysis. Table 1 summarizes the datasets used for training and validation.

All audio clips were downsampled to 16 kHz for training. Audio clips longer than 6 seconds were removed from the training and validation sets because they can exhaust available memory. We also removed audio samples of low quality that are labeled with more than one negative vote. After this preprocessing, we prepared three sets of training and validation data with different sample sizes. The first set contains 1k training samples and 200 validation samples per language, the second set 2k training samples and 400 validation samples per language, and the third set all the training samples available in the dataset. The samples were randomly selected from the preprocessed datasets, following the dataset splits provided by CommonVoice. For the evaluation of the model performance in supervised settings, i.e., where the languages used in the test are learned during the training, we picked 100 samples per language (i.e., 700 samples in total) from the test split.

For testing the universality of the models, in addition to the in-domain evaluation of the performance in the trained languages, we picked four low-resource languages from CommonVoice that are typologically diverse, geographically dis-

tant, and genetically unrelated: Luganda (< Bantu; Uganda), Upper Sorbian (< Indo-European; Saxony, Germany), Hakha Chin (< Sino-Tibetan; Chin State, Myanmar), and Tatar (< Turkic; Tatarstan, Russia). For annotation, we hired two language technology graduate students of their 20s who had had training in phonetics and IPA transcription, and we had them transcribe the test audio clips manually as mock fieldworkers. Their first language was Japanese. To prevent the annotator from guessing a language, they were not told what and how many languages were included in the samples. These clips were randomly chosen from these four languages, and the annotator transcribed the assigned audio clips until we had 100 annotated samples. We used the transcriptions by one of the two annotators as the gold labels in the evaluation, and those by the other were used to measure the inter-annotator agreement (IAA) between the two annotators with the metrics used in this study (see Section 4.2).

### 3.2. Model

The pretrained model is wav2vec2-large-xlsr-53 [5] which was trained on 56k hours of speech data with 53 languages from CommonVoice, Multilingual LibriSpeech, and BABEL. Our objective is to fine-tune the model to predict the IPA string of audio input. We implement the fine-tuned model with `Wav2Vec2ForCTC` provided in the `transformers` library to train it with the Connectionist Temporal Classification (CTC) loss [18]. Since our goal is to train a model applicable to unseen languages, we did not use encoder-decoder models that are not language-agnostic such as Whisper [7].

## 4. Experiments

### 4.1. Setup

We used `Wav2Vec2CTCTokenizer` as the tokenizer and included in the vocabulary the full list of the IPA characters and their possible combinations such as multi-letter phones and co-articulated consonants. The full IPA list, obtained from PanPhon [19], consists of 6,487 phones.

We set the CTC loss reduction (`ctc_loss_reduction`) to mean, the learning rate to $3.0 \times 10^{-4}$, the warmup steps to 500, the number of training epochs to 30; other numerical hyperparameters used the default values defined in the configuration in the `transformers` library as of version 4.26.0. We used 150 training epochs for the extremely low-resource setting with 100 samples per language and 5 epochs for the setting with full training samples, so that the models would learn to output IPA sequences well enough while keeping them from overfitting. The feature extractor of the pretrained model is frozen before fine-tuning.

In the default setting (1k training samples per language), the

---

[3] https://huggingface.co/datasets/mozilla-foundation/common_voice_11_0

[4] https://github.com/dmort27/epitran

Table 3: *The evaluation scores in the supervised setting for models with 1k, 2k, and full training samples.*

| Metric | # Train samples | Japanese | Polish | Maltese | Hungarian | Finnish | Greek | Tamil | Overall |
|---|---|---|---|---|---|---|---|---|---|
| PER (%) | 1k | 17.4 | 16.5 | 14.3 | 44.9 | 29.6 | 24.7 | 26.9 | 24.9 |
| | 2k | 12.8 | 12.3 | 10.4 | 43.1 | 28.3 | 23.7 | 25.8 | 22.4 |
| | full | 47.8 | 7.1 | 14.9 | 42.3 | 13.8 | 7.7 | 13.4 | 21.0 |
| PFER (%) | 1k | 4.8 | 4.5 | 5.2 | 11.7 | 8.2 | 8.4 | 6.2 | 7.0 |
| | 2k | **3.6** | 3.8 | **3.5** | **10.8** | 8.0 | 7.7 | 5.7 | 6.2 |
| | full | 8.9 | **2.5** | 4.9 | **10.8** | **5.0** | **4.1** | **3.6** | **5.7** |

Table 4: *Comparison of the scores in the zero-shot setting. The numbers in parentheses next to our models refer to the number of training samples per language. Human (IAA) is the IAA scores between the two human annotators.*

| Metric | Model | Luganda | Upper Sorbian | Hakha Chin | Tatar | Overall |
|---|---|---|---|---|---|---|
| PER (%) | Allosaurus | 104.1 | 93.9 | 79.4 | 89.6 | 91.8 |
| | Wav2Vec2Phoneme | 64.0 | 66.1 | 70.0 | 63.0 | 65.8 |
| | Ours (1k) | 74.0 | 68.6 | 73.0 | 67.7 | 70.8 |
| | Ours (2k) | 77.0 | 69.4 | 72.7 | 67.5 | 71.6 |
| | Ours (full) | 70.9 | 69.8 | 68.7 | 63.2 | 63.2 |
| | Human (IAA) | 52.9 | 52.5 | 55.3 | 52.7 | 53.3 |
| PFER (%) | Allosaurus | 46.1 | 36.3 | 36.3 | 30.1 | 34.2 |
| | Wav2Vec2Phoneme | 24.2 | 26.1 | **19.3** | 20.0 | 22.4 |
| | Ours (1k) | **20.8** | 24.0 | 21.3 | **18.8** | **21.2** |
| | Ours (2k) | 22.7 | 24.9 | 21.8 | 19.4 | 22.2 |
| | Ours (full) | 23.0 | **23.1** | 20.3 | **18.8** | 21.3 |
| | Human (IAA) | 19.3 | 22.1 | 17.8 | 19.1 | 19.6 |

training took ~4 hours in runtime with four GTX1080Ti GPUs, and the average power consumption was ~7.5kW.

### 4.2. Evaluation metrics

We compare our models with the two existing speech-to-IPA models: Allosaurus [14] and Wav2Vec2Phoneme [3]. We evaluate these models with two metrics: naïve phone error rate (PER) and PER with phonetic feature edit distance [19], which we call PFER (phone feature error rate) for short. PER is a modification of character error rate (CER) whose basic unit is a phone (which may consist of more than one Unicode character) instead of a character. However, a problem of using PER in this task is that PER ignores phonetic similarities among phones. For example, suppose that a model predicts [kʰæt] for a spoken word with the label [kæt], where they express the first consonant ([k] or [kʰ]) in a different manner. They only differ in one phonetic feature, namely aspiration, and share a similar acoustic impression; therefore, intuitively speaking, their distance should be smaller than that of two utterly unrelated phones like [k] and [a]. PER treats both pairs of phones equally. PFER, in contrast, considers acoustic similarities among phones by representing each phone as a collection of phonetic features. In our experiments, we used the implementation provided in Pan-Phon [19], which defines 24 phonetic features for each phone. PFER calculates the Hamming distance of features between two phones so that one feature mismatch has a cost of 1/24.[5] Other string operations (insertion and substitution) cost 1. For this reason, we put more emphasis on the PFER-based comparison than PER in this study because PFER is more representative of the transcription accuracy.

## 5. Results

This section reports the results of the performance of our models and compares them to the two existing speech-to-IPA models, Allosaurus and Wav2Vec2Phoneme, as well as the human IAA.

### 5.1. In-domain evaluation

Table 3 reports the performance of automatic IPA transcription for the languages the models were trained with. In both PER and PFER, our model performed better when the training size is larger. In the overall scores, the model with the maximum training size (50,348 samples, or ~60 hours, in total) achieved 5.7% PFER, while the models with 1k and 2k training samples per language scored 7.0% and 6.2% PFER, respectively. These results follow the general assumption in machine learning that training with more samples gives better performance.

### 5.2. Zero-shot evaluation

Table 4 compares the scores of the performance in transcribing languages unseen during the training phase. Overall, our model with 1k training samples per language (21.2% PFER) performed the best among the tested models. In particular, it outperformed both of Allosaurus (34.3% PFER) and Wav2Vec2Phoneme (22.4% PFER). Interestingly, our model with 1k training samples per language performed better on average than that with 2k samples and full samples despite its smaller size. Given that the overall human IAA rate is 19.6% PFER, our models and Wav2Vec2Phoneme have nearly reached performance comparable with human annotators. For reference, Table 5 compares the actual output by Allosaurus,

---

[5]The documentation of PanPhon at the time of writing notes that the cost of a feature mismatch is 1/22, but the implementation gives 1/24 because their feature table contains 24 features.

Table 5: *Output by Allosaurus, Wav2Vec2Phoneme, and our model (1k). The sample is extracted from the Luganda validation set in CommonVoice.*

| Sentence | Omukama mulungi obudde bwonna. |
|---|---|
| Allosaurus | omuək<sup>h</sup>amamuərondʒioʁuədʲeɓuəɔenːʌ |
| Wav2Vec2Phoneme | omukamamudundʒipovudevonna |
| Ours (1k) | omukamamurundʒipoputdɛvonna |

Wav2Vec2Phoneme, and our model (1k).

# 6. Discussion

The above results show that the performance of Wav2Vec2Phoneme and our models is comparable to human annotators, while still lagging behind by approximately 2 to 4% in PFER. The better performance of our model (1k) than Wav2Vec2Phoneme is particularly remarkable because ours was only trained on 7k samples (~9.7 hours) from seven languages, while Wav2Vec2Phoneme uses at least 100 hours of audio from at least 26 languages from CommonVoice. In addition, these results suggest that transfer learning for the speech-to-IPA task reaches a performance plateau with a rather small training dataset, and increasing the number of samples would not contribute to improvement. In the following ablation studies, we discuss what are the factors that affect performance.

## 6.1. Ablation studies

We considered four parameters in this ablation: training samples per language (Size), additional data from Forvo[6] (+/−Forvo), quality filtering (+/−QF), and linguistic diversity in the training data (# Languages). We used the same evaluation methods as in the zero-shot evaluation in Section 5.2. We prepared an additional small dataset with 353 samples from six languages (Adyghe, Arabic, Burmese, Icelandic, Xhosa, and Zulu) retrieved from Forvo to add phonetic and linguistic diversity to the training data. However, a numerical comparison in Table 6 shows no benefit from adding the Forvo data. It may be that the additional data size is too small to affect the performance, and further investigation is called for.

Within the setting with filtering out poor-quality audio (+QF in Table 6), the extremely low-resource setting with only 100 samples per language (i.e., 700 samples in total) performed the worst; however, the results with 1k samples were better than that with 2k samples. This implies that increasing the amount of training data does not promise a performance improvement and that it might hit a plateau before 1k samples per language. Table 6 also shows that removing audio files of poor quality gave us better results with 1k samples per language but did worse with 2k samples per language. This suggests that removing audio files of poor quality does not necessarily promise an improvement, either.

Table 7 specifically compares the effect of linguistic diversity in the training data. We can see that having more diversity provides better results. This is expected since having fewer languages in the training set means that predictions will be strictly limited to the phonetic inventory of those languages.

Table 6: *Ablation studies on quality filtering and additional data. The unit is PFER (%). "+/−Forvo" means whether the Forvo dataset was included in training, and "+/−QF" whether low-quality audio samples were removed. The Forvo dataset was not applied to the low-resource setting (Size=100).*

| | +Forvo | | −Forvo | |
|---|---|---|---|---|
| Size | +QF | −QF | +QF | −QF |
| 100 | — | — | 24.3 | 22.6 |
| 1k | 22.6 | 21.4 | 21.2 | 21.2 |
| 2k | 21.8 | 22.5 | 22.2 | 22.9 |

Table 7: *Ablation studies on the effect of linguistic diversity in the training data. The unit is PFER (%). All models were trained without audio of bad quality and the Forvo dataset.*

| # Languages | PFER (%) |
|---|---|
| 1 (ja) | 28.3 |
| 3 (ja, pl, mt) | 24.1 |
| 7 (all) | 21.2 |

## 6.2. Limitations

Last but not least, we mention several limitations that this study has faced but can be improved in future research. First, we could only use the seven languages that can be automatically transcribed from graphemes to phones relatively easily and accurately (Section 3.1). The languages used to train the model are still too few to reflect the actual linguistic diversity of the world. Second, it is not a perfect solution to rely on rule-based G2P tools to generate IPA transcriptions to be used as the labels, because our objective is to generate IPA transcription as accurately as human transcription. Third, due to the difficulty of hiring well-trained annotators for IPA transcription, our test data for the zero-shot evaluation only contains 100 samples in total. We can overcome these three limitations by developing larger high-quality datasets with speech and phonetic transcriptions. Fourth, our models do not consider tones as phonetic features. Extensions to include lexical and grammatical tones can be done by fine-tuning the model to the phonology of a specific tonal language; or, more generally, it might be possible to incorporate tones in the system by adding a layer or token that implicitly identifies whether the language is tonal or not.

# 7. Conclusion

This study introduced a new universal speech-to-IPA model trained with only 7k samples from seven languages. We showed that our model with 1k training samples per language performs better than existing speech-to-IPA models trained on larger datasets from at least 26 languages. Our model achieved 21.2% in PFER, which was almost comparable to the human IAA score. In our settings, it was observed that broader linguistic diversity in the training data gives more accurate IPA transcription. Our results also suggested the importance of using clean phonetic transcription in the training dataset.

# 8. Acknowledgments

---

[6]https://forvo.com

# 9. References

[1] P. Austin and J. Sallabank, Eds., *The Cambridge Handbook of Endangered Languages*, ser. Cambridge Handbooks in Language and Linguistics. Cambridge University Press, 2011. [Online]. Available: https://doi.org/10.1017/CBO9780511975981

[2] N. Thieberger, "LD&C possibilities for the next decade," *Language Documentation & Conservation*, vol. 11, pp. 1–4, 2017.

[3] Q. Xu, A. Baevski, and M. Auli, "Simple and effective zero-shot cross-lingual phoneme recognition," in *Proceedings of Interspeech*, 2022, pp. 2113–2117. [Online]. Available: https://doi.org/10.21437/Interspeech.2022-60

[4] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020, pp. 12 449–12 460. [Online]. Available: https://arxiv.org/abs/2006.11477

[5] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," in *Proceedings of Interspeech*, 2021, pp. 2426–2430. [Online]. Available: https://doi.org/10.21437/Interspeech.2021-329

[6] B. Li, R. Pang, Y. Zhang, T. N. Sainath, T. Strohman, P. Haghani, Y. Zhu, B. Farris, N. Gaur, and M. Prasad, "Massively multilingual ASR: A lifelong learning solution," in *Proceedings of the 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 6397–6401. [Online]. Available: https://doi.org/10.1109/ICASSP43922.2022.9746594

[7] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," 2022. [Online]. Available: https://arxiv.org/abs/2212.04356

[8] Y. Zhang, W. Han, J. Qin, Y. Wang, A. Bapna, Z. Chen, N. Chen, B. Li, V. Axelrod, G. Wang, Z. Meng, K. Hu, A. Rosenberg, R. Prabhavalkar, D. S. Park, P. Haghani, J. Riesa, G. Perng, H. Soltau, T. Strohman, B. Ramabhadran, T. Sainath, P. Moreno, C.-C. Chiu, J. Schalkwyk, F. Beaufays, and Y. Wu, "Google USM: Scaling automatic speech recognition beyond 100 languages," 2023. [Online]. Available: https://arxiv.org/abs/2303.01037

[9] T. Schultz and A. Waibel, "Fast bootstrapping of lvcsr systems with multilingual phoneme sets," in *Proceedings of 5th European Conference on Speech Communication and Technology (EUROSPEECH '97)*, vol. 1, September 1997, pp. 371 – 373. [Online]. Available: https://www.isca-speech.org/archive_v0/eurospeech_1997/e97_0371.html

[10] P. Cohen, S. Dharanipragada, J. Gros, M. Monkowski, C. Neti, S. Roukos, and T. Ward, "Towards a universal speech recognizer for multiple languages," in *Proceedings of the 1997 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, 1997, pp. 591–598. [Online]. Available: https://doi.org/10.1109/ASRU.1997.659140

[11] T. Schultz and A. Waibel, "Language-independent and language-adaptive acoustic modeling for speech recognition," *Speech Communication*, vol. 35, no. 1, pp. 31–51, 2001. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0167639300000947

[12] S. Dalmia, R. Sanabria, F. Metze, and A. W. Black, "Sequence-based multi-lingual low resource speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE Press, 2018, p. 4909–4913. [Online]. Available: https://doi.org/10.1109/ICASSP.2018.8461802

[13] A. Michaud, O. Adams, T. A. Cohn, G. Neubig, and S. Guillaume, "Integrating automatic transcription into the language documentation workflow: Experiments with Na data and the Persephone toolkit," *Language Documentation & Conservation*, vol. 12, pp. 481–513, 2018. [Online]. Available: http://hdl.handle.net/10125/24793

[14] X. Li, S. Dalmia, J. Li, M. Lee, P. Littell, J. Yao, A. Anastasopoulos, D. R. Mortensen, G. Neubig, A. W. Black, and M. Florian, "Universal phone recognition with a multilingual allophone system," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 8249–8253. [Online]. Available: https://doi.org/10.1109/ICASSP40776.2020.9054362

[15] H. Gao, J. Ni, Y. Zhang, K. Qian, S. Chang, and M. Hasegawa-Johnson, "Zero-shot cross-lingual phonetic recognition with external language embedding," in *Proceedings of Interspeech*, 2021, pp. 1304–1308. [Online]. Available: https://doi.org/10.21437/Interspeech.2021-1843

[16] B. Michailovsky, M. Mazaudon, A. Michaud, S. Guillaume, A. François, and E. Adamou, "Documenting and researching endangered languages: The Pangloss Collection," *Language Documentation & Conservation*, vol. 8, pp. 119–135, 2014. [Online]. Available: http://hdl.handle.net/10125/4621

[17] D. R. Mortensen, S. Dalmia, and P. Littell, "Epitran: Precision G2P for many languages," in *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 2018. [Online]. Available: https://aclanthology.org/C16-1328

[18] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd International Conference on Machine Learning*, 2006, p. 369–376. [Online]. Available: https://doi.org/10.1145/1143844.1143891

[19] D. R. Mortensen, P. Littell, A. Bharadwaj, K. Goyal, C. Dyer, and L. S. Levin, "PanPhon: A resource for mapping IPA segments to articulatory feature vectors," in *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 2016, pp. 3475–3484. [Online]. Available: https://aclanthology.org/C16-1328