
How to make someone speak a language that they don't know.

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We present a simple idea that allows to record a speaker in a given language and
2 synthesize their voice in other languages that they may not even know. These
3 techniques open a wide range of potential applications such as cross-language
4 communication, language learning or automatic video dubbing. We call this
5 general problem *multi-language speaker-conditioned speech synthesis* and we
6 present a simple but strong baseline for it.

7 Our model architecture is similar to the encoder-decoder Char2Wav model [Sotelo
8 et al., 2017] or Tacotron [Shen et al., 2017]. The main difference is that, instead
9 of conditioning on characters or phonemes that are specific to a given language,
10 we condition on a shared phonetic representation that is universal to all languages
11 [Meier, 2016]. This cross-language phonetic representation of text allows to
12 synthesize speech in any language while preserving the vocal characteristics of the
13 original speaker. Furthermore, we show that fine-tuning the weights of our model
14 allows us to extend our results to speakers outside of the training dataset.

15 1 Introduction

16 The goal of this work is to create a text-to-speech system able to generate audio in multiple languages
17 for any given speaker. We further impose two requirements. First, the model should be able to copy
18 the voice of an out-of-dataset speaker given only very limited data. Second, the model should be able
19 to generate audio in any language, even when trained on a single-language speaker.

20 Such a system, paired with a word translation system, would enable anyone to speak in any language.
21 It could be used by travelers to help them communicate in foreign countries or by movie producers to
22 dub movies while keeping the original voices of their actors.

23 Our approach is to build a model able to generate speech in multiple languages. The model is trained
24 with multiple speakers to let the model be aware of the variations between speakers and also to
25 disentangle speech content from speaker identity. Once the model is trained, we bias the generation
26 process so that it sounds like a specific speaker. This speaker doesn't have to be in the training data.

27 2 Related Work

28 Our work builds upon recent developments in neural network based speech synthesis [Sotelo et al.,
29 2017, Ping et al., 2017, Shen et al., 2017, Van Den Oord et al., 2016]. Specifically, our model
30 architecture closely resembles attention-based speech synthesis models, which map sequences of
31 phonemes or characters to intermediate audio representations e.g. vocoder [Morise et al., 2016][Sotelo
32 et al., 2017] or spectrogram [Shen et al., 2017][Ping et al., 2017]. The representation is post-processed
33 via either signal processing based methods e.g. Griffin-Lim spectrogram inversion [Griffin and Lim,

34 1984], World/Straight vocoder [Morise et al., 2016] or neural vocoders [Sotelo et al., 2017][Shen et al.,
 35 2017][Ping et al., 2017] to get raw audio. Our model is very closely related to neural multi-language
 36 multi-speaker parametric speech synthesis model described by Li and Zen [2016]. Our approach
 37 enables speaker style transfer across languages present in the training dataset, which is different
 38 from generalising to new unseen languages with small amount of data. Also Li and Zen [2016] uses
 39 union of language-dependent linguistic feature set to represent text input to their model. As opposed
 40 to them, we use International Phonetic Alphabet (IPA) [Meier, 2016] to map text input across all
 41 languages to a universal representation.

42 Our model is able to accomplish zero-shot accent transfer, which is very similar to zero-shot machine
 43 translation, done by grounding the input from different languages to a common neural representation
 44 space, followed by decoding in the audio space [Johnson et al., 2016].

45 3 International Phonetic Alphabet (IPA)

46 In this work we use the International Phonetic Alphabet (IPA) [Meier, 2016] to represent the text
 47 information. The IPA is designed to represent only those qualities of speech that are part of oral
 48 language: phones, phonemes, intonation and the separation of words and syllables. Note that the
 49 IPA is language independent. It thus provides a common representation for the text input across
 50 several languages to be conditioned upon for speech generation. In the following, we will show that
 51 our model is able to generalize the phonetic information learned across multiple languages. This
 52 reduces the complexity of text conditioning and improves data efficiency when we train the model
 53 with multilingual data.

54 4 Model

55 At the core of our Text-To-Speech system lies an attention-based encoder-decoder architecture similar
 56 to [Sotelo et al., 2017] [Shen et al., 2017] [Bahdanau et al., 2014]. More specifically, we use a
 57 bidirectional recurrent conditioning encoder, a recurrent decoder and a location based attention
 58 mechanism very similar to the one used in [Graves, 2013]. Discrete conditioning information i.e. IPA
 59 sequences and speaker information, are modeled with randomly initialized embedding layers that are
 60 trained together with the rest of the model.

61 The training data consists of audio-transcript pairs. The transcript is translated into its IPA equivalent
 62 before being fed to the model and the audio is transformed into an intermediate representation (e.g.
 63 WORLD vocoder parameters or spectrogram). Each speaker within the training dataset only speaks
 64 a single language. However, at synthesis time, we are able to take any combination of speaker and
 65 language, and produce natural sounding speech in the voice of the speaker and in the accent matching
 66 that of the language.

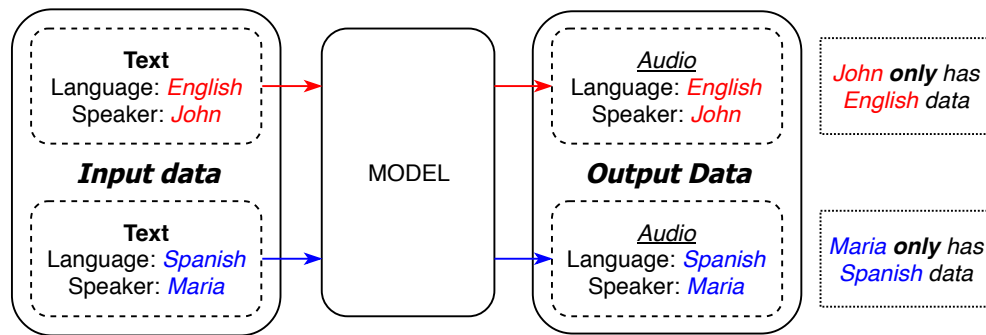


Figure 1: Training process. The model is trained to map text to audio based on (*transcript-audio*) pairs. Maria is a speaker for whom the model only sees Spanish data, John is a speaker for whom the model only sees English data

67 During inference, the model is able to generate Spanish audio in the style of an English speaker and
 68 vice-versa.



(a) Generating Spanish audio in the style of the English speaker (b) Generating English audio in the style of the Spanish speaker

Figure 2: Inference process. The model is able to interpolate and generate any (*speaker, language*) combination.

69 The method presented above only allows us to generate speech for speakers that are present in the
 70 training dataset. In order to generate multilingual speech for a new speaker, we fine-tune the trained
 71 model on a small amount of new speaker’s data, which can be as little as merely 300 seconds of
 72 utterances. It is important to preventing catastrophic forgetting during fine-tuning. We carried out
 73 a careful and exhaustive search hyperparameter search to obtain robust performance across many
 74 speakers. Crucially, we apply a smaller learning rate to the encoder and decoder parts of the models,
 75 and a higher one for the speaker embedding. This improved speaker fidelity considerably.

76 After fine-tuning, the model is able to generate any text in any language¹ with the new speaker’s
 77 vocal identity.

78 5 Experiments

79 We conduct experiments on our models trained in two distinct settings. First, we train our model
 80 with data in two languages (Bilingual Model). Second, we train our model with data in six languages
 81 (Multilingual Model).

82 For these experiments, we used several datasets. We used an internal English dataset composed of
 83 approximately 20000 speakers, with about 10 utterances per speaker. We also used the TIMIT dataset
 84 [Garofolo et al., 1993] and DIMEx100 [Pineda, 2009]. DIMEx100 is a Spanish dataset composed of
 85 100 Spanish native speakers, with about 60 2-seconds utterances per speaker.

86 For all the experiments we provide audio samples² rather than an exhaustive quantitative analysis.

87 5.1 Bilingual Model

88 We use a bilingual dataset (English and Spanish) composed of TIMIT and DIMEx100 to train this
 89 model. We concatenate to this data from the speaker we want to clone. We are able to generate
 90 bilingual speech for both English and Spanish speakers within the dataset. We also fine-tune the
 91 trained model on new speaker’s data to generate bilingual speech for the new speaker. We only need
 92 data from the new speaker in a **single** language.

93 5.2 Multilingual model

94 For this experiment we train the model with our internal multi-speaker dataset to which we add data
 95 from 5 single-speaker audiobooks from the CSS10 dataset [Kyubyong Park, 2018], each in a distinct
 96 language. We show that the model is able to generate in any language for any speaker in the dataset.
 97 The model also shows robust performance on new, out-of-sample speakers after the fine-tuning step
 98 (see figure 3).

99 References

100 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
 101 learning to align and translate. *CoRR*, abs/1409.0473, 2014. URL <http://arxiv.org/abs/1409.0473>.

¹Included in the training dataset.

²<https://everyone-speaks-every-language.github.io/>

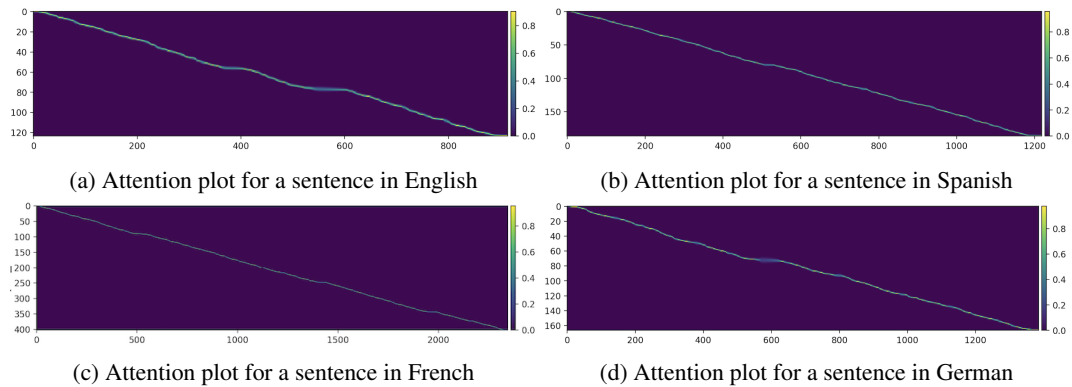


Figure 3: Attention plots for the same sentence in English, Spanish, French and German (same speaker)

- 103 J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, and N. L. Dahlgren. Darpa timit
104 acoustic phonetic continuous speech corpus cdrom, 1993.
- 105 Alex Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013.
106 URL <http://arxiv.org/abs/1308.0850>.
- 107 Daniel Griffin and Jae Lim. Signal estimation from modified short-time fourier transform. *IEEE*
108 *Transactions on Acoustics, Speech, and Signal Processing*, 32(2):236–243, 1984.
- 109 Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil
110 Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. Google’s multilingual neural
111 machine translation system: enabling zero-shot translation. *arXiv preprint arXiv:1611.04558*,
112 2016.
- 113 Tommy Mulc Kyubyong Park. Csx10: A collection of single speaker speech datasets for 10 languages.
114 <https://github.com/Kyubyong/csx10>, 2018.
- 115 Bo Li and Heiga Zen. Multi-language multi-speaker acoustic modeling for lstm-rnn based statistical
116 parametric speech synthesis. In *INTERSPEECH*, pages 2468–2472, 2016.
- 117 Paul Meier. *International Phonetic Alphabet (IPA) Charts*. Paul Meier Dialect Services, Lawrence,
118 KS, 2016. URL <http://www.paulmeier.com/ipacharts/>. bibtex: meier_international_2016.
- 119 Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. World: a vocoder-based high-quality speech
120 synthesis system for real-time applications. *IEICE TRANSACTIONS on Information and Systems*,
121 99(7):1877–1884, 2016.
- 122 Hayde Cuetara Javier Galescu Lucian Juárez Janet Llisterra Joaquim Pérez Patricia Villaseñor-
123 Pineda Luis Pineda, Luis Castellanos. The corpus dimex100: Transcription and evaluation.
124 *Language Resources and Evaluation*. 44. 347-370. 10.1007/s10579-009-9109-9, 2009.
- 125 Wei Ping, Kainan Peng, Andrew Gibiansky, Sercan O Arik, Ajay Kannan, Sharan Narang, Jonathan
126 Raiman, and John Miller. Deep voice 3: 2000-speaker neural text-to-speech. *arXiv preprint*
127 *arXiv:1710.07654*, 2017.
- 128 Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng
129 Chen, Yu Zhang, Yuxuan Wang, R. J. Skerry-Ryan, Rif A. Saurous, Yannis Agiomyrgiannakis,
130 and Yonghui Wu. Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions.
131 *CoRR*, abs/1712.05884, 2017. URL <http://arxiv.org/abs/1712.05884>.
- 132 Jose Sotelo, Soroush Mehri, Kundan Kumar, Joao Felipe Santos, Kyle Kastner, Aaron Courville, and
133 Yoshua Bengio. Char2wav: End-to-end speech synthesis. 2017.
- 134 Aäron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,
135 Nal Kalchbrenner, Andrew W Senior, and Koray Kavukcuoglu. Wavenet: A generative model for
136 raw audio. In *SSW*, page 125, 2016.