# PROSPECTS OF ARTICULATORY TEXT-TO-SPEECH SYNTHESIS

*Simon Stone[1], Atilla Azgin[2], Sabrina Mänz[2], Peter Birkholz[1]*
*[1]Technische Universität Dresden, Germany*
*[2]Aristech GmbH, Heidelberg, Germany*
*simon.stone@tu-dresden.de*

**Abstract:** Text-to-Speech (TTS) synthesis has come a long way from the first (English) system by Matsui et al. [1], which used a cascade of analyses and models to get from a text representation to the acoustic speech signal, to the currently popular end-to-end systems (e.g., Tacotron [2]), which use a single large model (usually a recurrent neural network) to map input character sequences directly to the corresponding acoustic waveform. However, all current TTS systems use pre-recorded speech snippets (at least during the training process), which means that they are limited to the available units (in unit-selection systems) or at least to the voice and speaking style of the training material (in end-to-end systems). Any para-linguistic expressiveness, e.g., emotion, can only be "applied" after the "neutral" signal is synthesized using morphing techniques, which introduce audible signal-processing artifacts and can lower the perceived naturalness. In this paper, we therefore present our ongoing work on a TTS pipeline that uses an articulatory-based synthesis backend instead of a signal-based synthesizer. We use the VocalTractLab*'s [3] vocal tract, articulatory, and acoustic models to perform the synthesis. The input modality for the articulatory control is the so-called gestural score [4]: a sequence of parallel articulatory targets that are approximated at adjustable speed by the individual articulators (see Figure 1). While this kind of control over the vocal tract model is very powerful, flexible, and expressive, offering arguably the same degrees of freedom as the human vocal tract itself, it is also very time consuming to arrange the score for even a basic utterance and requires expert knowledge of the phonetic and articulatory intricacies. Therefore, articulatory synthesis has so far been confined to a niche in speech research despite its many advantages over unit-selection or waveform synthesis.

In this paper, we outline our efforts to bring articulatory synthesis closer to the mainstream by introducing a preprocessing pipeline that generates a gestural score (and thus, by extension, the speech signal) from plain, orthographic text. The first stage in an articulatory TTS pipeline is similar to, e.g., unit selection: The transliteration of the orthographic text into a syllabified, phonetic representation (grapheme-to-phoneme conversion). Currently, this stage uses a proprietary web-service by our partner Aristech GmbH using a hand-transcribed lexicon and a Finite-State-Transducer to handle out-of-vocabulary words. After this stage, the pipeline diverges from the usual procedure, because articulatory synthesis requires very low-level information: Every articulatory gesture requires a target position, a duration, and a time constant (controlling the speed of the articulatory target approximation). The sequence of target positions can be derived from the phonetic transcription, since every phoneme corresponds uniquely to a vocal tract configuration. For the timing information, the acoustic duration of each syllable is determined using a Deep Neural Network and a combination of phonetic and linguistic features (details in the full paper). However, the mapping was trained using acoustic sound durations, because no training data to predict the gestural durations directly is currently available. Therefore, these durations are only used to initialize the gestures and their true timing is found in an iterative way to match the simulated articulatory landmarks with the given acoustic landmarks (phone boundaries). Finally, the gestural score is generated and used to synthesize

---

* http://www.vocaltractlab.de

the speech signal. Based on this neutral gestural score, many manipulations can be introduced to add para-linguistic information (e.g., changing the pitch offset and range, the articulatory precision, the phonation type, the subglottal pressure, and many more). In contrast to morphing-based signal processing techniques to turn a neutral speech *signal* into an expressive speech signal, manipulations at the *articulatory* level are much more realistic sounding and can focus on individual features without unintentional side effects. Some audio examples of such manipulations are available at http://www.vocaltractlab.de/index.php?page=birkholz-supplements.
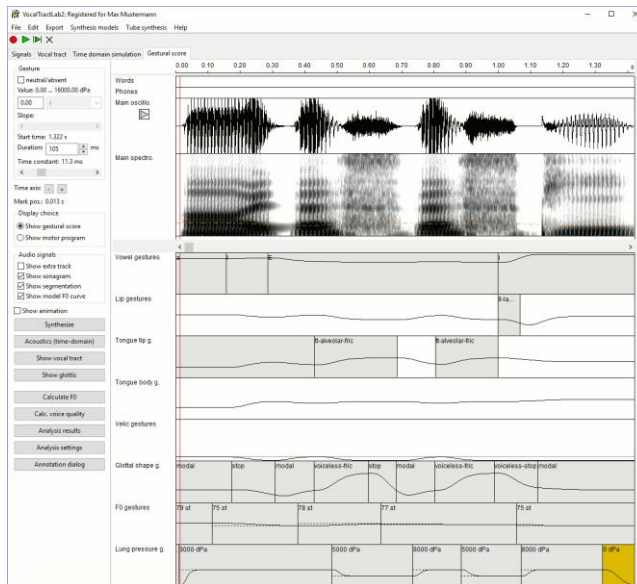


**Figure 1:** Gestural score for the utterance "ISSP" (/aI ?Es ?Es pi:/).

**References**

[1] Umeda, Noriko, E. Matsui, Torazo Suzuki, and Hiroshi Omura. "Synthesis of fairy tales using an analog vocal tract." In *Proceedings of the 6th International Congress on Acoustics*, pp. B159-162. 1968.

[2] Wang, Yuxuan, R. J. Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Zongheng Yang Jaitly, Ying Xiao et al. "Tacotron: Towards End-to-End Speech Synthesis." In *Proc. of the Interspeech*, Stockholm, Sweden, pp. 4006-4010. 2017.

[3] Birkholz, Peter. "Modeling consonant-vowel coarticulation for articulatory speech synthesis." *PloS one* 8, no. 4 (2013): e60603.

[4] Browman, Catherine P., and Louis Goldstein. "Articulatory phonology: An overview." *Phonetica* 49, no. 3-4 (1992): 155-180.

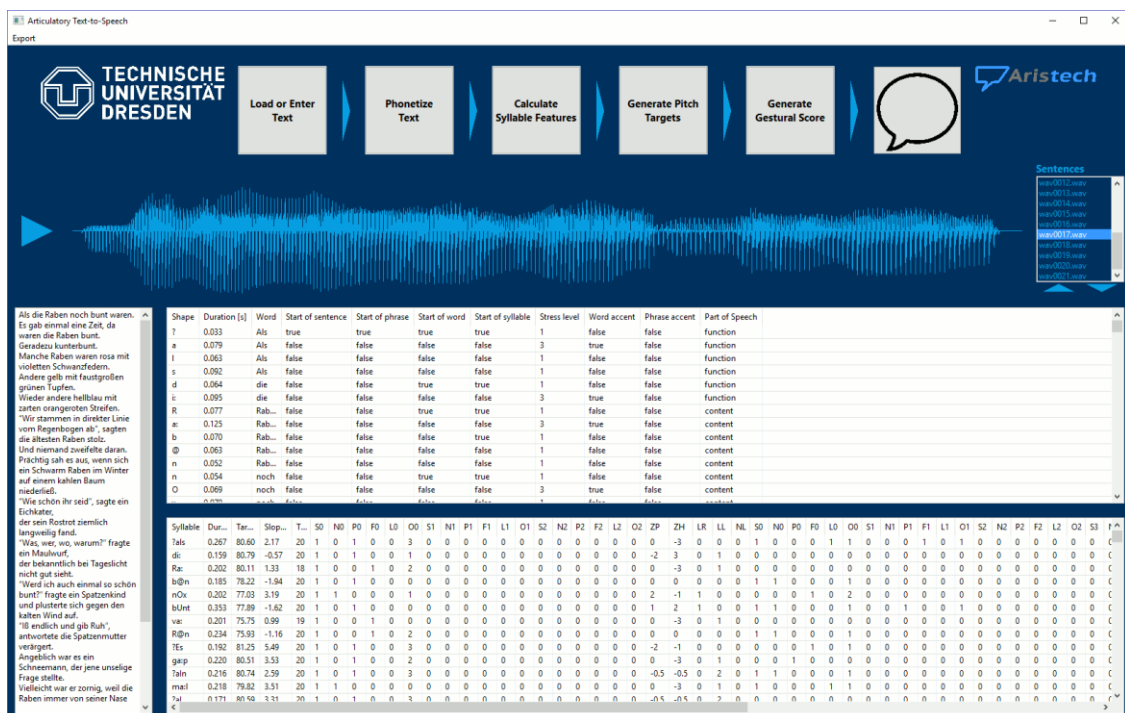**Figure 2:** Interface for the Articulatory Text-to-Speech Synthesis.