

Final Report of the NKFIH/FK-17-124584 project
‘Silent Speech Interface based on articulatory movements’
Sep 1, 2017 – Feb 28, 2022

1. Introduction

During the last several years, there has been significant interest in the articulatory-to-acoustic conversion research field, which is often referred as “Silent Speech Interfaces” (SSI) [1]. This has the main idea of recording the articulatory movement, and automatically generating speech from the movement information, while the original subject is not producing any sound. Such an SSI system can be highly useful for the speaking impaired (e.g., after laryngectomy), and for scenarios where regular speech is not feasible but information should be transmitted from the speaker (e.g., extremely noisy environments; military applications).

During the FK-17 OTKA project, we published 27 conference papers (the majority at top-ranked international conferences like Interspeech, IJCNN, Speech Synthesis Workshop), 7 international journal papers (mostly with Q1/Q2 ranking, sum impact factor: 14.919; and two more manuscripts are still under revision). Nine BSc, 13 MSc and 5 PhD students were involved as part of their project laboratory, thesis writing or individual research project. We had significant discussions with international researchers at renowned conferences (e.g., Interspeech 2017, 2018, 2019, 2020, 2021, ISSP 2021, SSW11, AICV, SPECOM and DAGA 2021) about methods that they / we use, and also about potential cooperation possibilities (e.g., with Dr. Michael Pucher from ÖAW, Austria, a joint Austrian-Hungarian cooperation grant is foreseen in the topic of Voicebanks and Speech Interfaces).

With the cooperation of Prof. Bruce Denby (Sorbonne Université, France), who was the author of the first Silent Speech Interface related paper in 2010 [1], and Dr. Michael Wand (IDSIA, Switzerland), who is expert in advanced deep learning [2], we proposed a special issue of the MDPI Sensors journal, entitled "Future Speech Interfaces with Sensors and Machine Intelligence", https://www.mdpi.com/journal/sensors/special_issues/FSL-SMI. We contacted a large number of international colleagues in this field, and envisage receiving many significant manuscripts related to SSI. The call is still open until May 2022; currently there are four published papers, and two other submissions are under review or further processing.

Within the FK-17 and PD-18 projects, we mostly proposed methods which are applying Ultrasound Tongue Imaging (UTI) while contributing to the SSI field with. Originally, we expected to use Electromagnetic Articulography (EMA) as well, but after the grant submission, Prof Jun Wang's research lab (U.Texas at Austin, USA) achieved significant research results in EMA-to-speech generation and recognition [3]. Therefore, in the current project, instead of EMA, we focused on ultrasound, lip video and MRI of the vocal tract.

1.1 Key questions, goals of the project

The key goals of the project were to 1) thoroughly analyze the articulatory phone recognition performance using the optimal combination of different articulatory tracking methods, 2) enhance spectral filtering of vocoding using articulatory data, 3) test and improve recognition-and-synthesis and direct synthesis in the field of silent speech interfaces. Finally, we planned to create a prototype silent speech interface that is using articulatory data as input and generates speech as output.

2. Methods, experiments and results

The project was divided into four work packages (WP), according to the main goals, and the results are summarized separately for each WP. Within the project, there was a strong cooperation between several universities: Budapest University of Technology and Economics (the PI, and his BSc, MSc and PhD students), University of Szeged (Gábor Gosztolya, László Tóth, Tamás Grósz), ELTE & MTA-ELTE Lingual Articulation Research Group (Alexandra Markó). Besides, during the four and half years, several other Hungarian and international researchers or PhD

students joined for shorter periods (e.g., Csaba Zainkó and Géza Németh at BME, Amin Shandiz at SZTE, Dagoberto Porras from IUS, Colombia, etc.).

2.1 WP1: Articulatory movement based phone recognition

First, after the project start in 2017, we recorded a database with five male and four female Hungarian subjects for the later experiments (200 sentences, roughly 20 minutes data from each of them). Besides, we recorded a female speaker at five different sessions. The recordings involved the acquisition of tongue-ultrasound and lip video as articulatory data, and speech as well. The recordings were done at the facilities of ELTE, with the equipment of the MTA-ELTE Lingual Articulation Research Group, including the “Micro” ultrasound system. We recorded both regular speech (full sentences) and silent articulation, in midsagittal orientation. From the speech, the spectral features and F0 were extracted; whereas the ultrasound data was stored in raw scanline format (binary pixels, 64 scanlines, 842 points on each).

After the recordings were done, we conducted the following experiment related to phone recognition [4]. We recognized that the learning of speech recognition and speech synthesis targets (acoustic model states vs. vocoder parameters) are two closely related tasks over the same ultrasound tongue image input, so here we experimented with the multi-task training of deep neural networks (see Fig. 1.), which seeks to solve the two tasks simultaneously. Our results showed that the parallel learning of the two types of targets is indeed beneficial for both tasks. Moreover, we obtained further improvements by using multi-task training as a weight initialization step before task-specific training. Overall, we reported a relative error rate reduction of about 7% in both the speech recognition and the speech synthesis tasks [4].

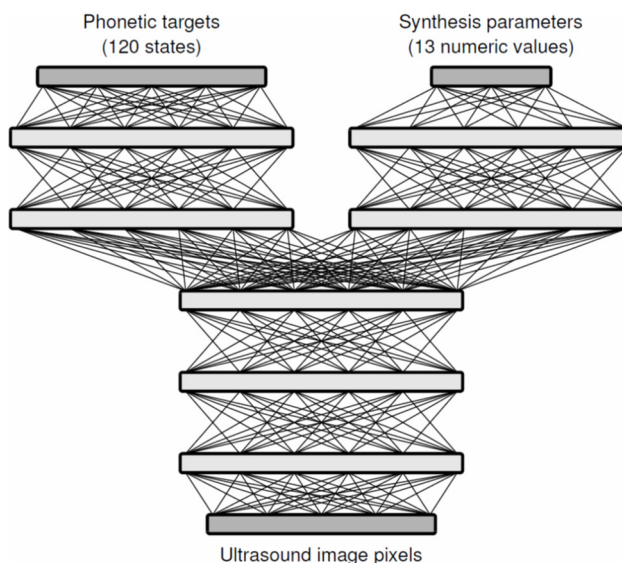


Fig. 1: Structure of the multi-task DNN, for the case of 3 shared and 2 task-specific layers. From [4]

2.2 WP2: Articulatory motivated spectral filtering in vocoding

During submission of the FK-17 project, Alexander Sepulveda (Industrial University of Santander, Bucaramanga, Colombia) signed a declaration of international cooperation. As a result of this, we hosted an international student at BME (Dagoberto Porras Plata, from the same university in Colombia) for professional internship between Oct – Dec 2018, who was involved in the acoustic-to-articulatory inversion (AAI) experiments of the project [5]. We implemented several different Deep Neural Networks (DNNs) to estimate the articulatory information from the acoustic signal. Ultrasound Tongue Imaging was used as the target articulatory information, and we tested two approaches: 1) the EigenTongue space and 2) the raw ultrasound image, and found that raw target data and a simple neural network with two hidden layers were more suitable for this inversion task. After that, we implemented several advanced DNNs (convolutional and recurrent neural networks), to estimate ultrasound images from the acoustic signal. From these experiments, a journal manuscript has been written, which is still under review [6]. Also, within

the AAI scenario, we tested real-time MRI of the vocal tract for the target of the DNNs [7]. As the result, LSTMs can achieve smooth generated MR images of the vocal tract, which are similar to the original MRI recordings. This way, we could compare two articulatory imaging methods (namely ultrasound and MRI) in the AAI field.

Next, in this WP, we experimented with novel vocoders for text-to-speech (TTS) synthesis and voice conversion (VC) with Mohammed Al-Radhi, being the PhD student of the PI. We proposed a continuous vocoder using continuous F0 (contF0) in combination with Maximum Voiced Frequency (MVF), and applied this for recurrent neural network based voice conversion [8]. The continuous vocoder was applied in DNN-based TTS and tested with both English and Arabic speech [9]. As an extension, Continuous Noise Masking (CNM) was proposed to overcome the issues of simple vocoders (e.g., buzziness) [10]. Within Statistical Voice Conversion (SVC), multiple features from the speech of two speakers (source and target) are converted, using DNNs [11]. We integrated into the SVC framework the continuous vocoder, by converting its contF0, maximum voiced frequency, and spectral features. The continuous vocoder, that we used earlier for UTI-to-F0 prediction [12], was further developed for speech synthesis and voice conversion [13]. We applied Continuous Wavelet Transform (CWT) to characterize and decompose speech features [14]. It can retain the fine spectral envelope and achieve high controllability of the structure closer to human auditory scale. Finally, a demo application 'conTTS' was created [15].

2.3 WP3: Automatic articulatory-to-acoustic mapping using deep learning methods

Although the main focus of this project was ultrasound as the articulatory acquisition technique, we compared several other types as well: ultrasound tongue imaging (UTI), lip video (LIP), and vocal tract Magnetic Resonance Imaging (MRI). The advantage of ultrasound is that the full tongue is visible with relatively good spatial and temporal resolution. Lip video, on the other hand, is much simpler to record; but contains significantly less information about the articulation. VT-RMI is significantly more complex and expensive to record, but can provide very detailed information about the full vocal tract.

Silent Speech Interface systems apply two different strategies to solve the articulatory-to-acoustic conversion task. The recognition-and-synthesis approach applies speech recognition techniques to map the articulatory data to a textual transcript, which is then converted to speech by a conventional text-to-speech system. The direct synthesis approach (i.e., WP3) seeks to convert the articulatory information directly to speech synthesis (vocoder) parameters. First, we used various neural network types for articulatory-to-speech synthesis: fully connected network [16] and convolutional and recurrent networks [17]. A pre-processing using a Deep Convolutional AutoEncoder was also studied. From these, an architecture based on a CNN and bidirectional LSTM layers has shown the best objective and subjective results [17]. We also tested more complex architectures, like 3D convolutional networks [18], and multi-task learning [4].

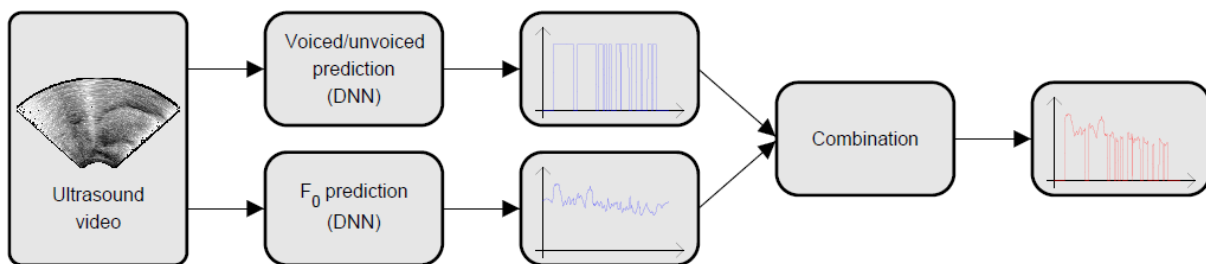


Fig. 2: Separate DNNs to estimate the voicing and the value of F0 from ultrasound input. From [19].

Next, we conducted experiments in ultrasound-to-F0 prediction: we compared a standard pulse-noise vocoder [19] (Fig. 2), several F0 estimation approaches [20] and also used a simple continuous F0 tracker which does not apply a strict voiced / unvoiced decision [12]. Continuous vocoder parameters (ContF0, Maximum Voiced Frequency and Mel-Generalized Cepstrum) were predicted using separate convolutional neural networks, with ultrasound as input. The methods were tested on four Hungarian speakers (2 males and 2 females), who were recorded in WP1. As further experiments, we proposed AutoEncoders for the representation of ultrasound tongue images [21]; and proposed voice activity detection from ultrasound images [22]. Also, we investigated the degree of session-dependency of standard

feed-forward DNN-based models for ultrasound-based SSI systems [23]. Besides examining the amount of training data required for speech synthesis parameter estimation, we also showed that DNN adaptation can be useful for handling session dependency [23]. Next, we applied a flow-based neural vocoder (WaveGlow) for ultrasound-to-speech synthesis [24]. Here, the training target was the 80-dimensional mel-spectrogram, which results in a finer detailed spectral representation than earlier methods [24]. Still, using UTI, we tested Generative Adversarial Networks (GAN) for generation of ultrasound images [25]. We compared ultrasound image representations (raw vs. wedge) as the input of SSI systems, of which a journal manuscript is still under review [26]. As a further experiment, we conducted recognition from ultrasound [27] and compared this with the direct synthesis approach and also with speech synthesis from text [28] (Fig. 3).

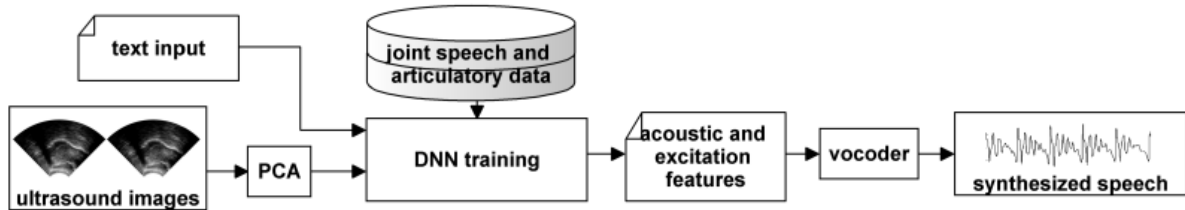


Fig. 3: Comparison of direct ultrasound-to-speech synthesis approach and speech synthesis from text. From [28].

As a second articulatory acquisition technique, we tested lip videos & mouth movement [29]. Inspired by earlier lip-to-speech studies, in our research we designed and implemented models that can generate spectral parameters of speech from lip videos. We used 1000 sentences from a male English speaker of the GRID audiovisual database [30], which contains video from the face of speakers, and synchronous speech. We tested two models that use convolutional and recurrent layers, of which the recurrent neural network was preferred. The results on lip processing might be useful for the FK-17 project, as recording the lip video is simple and cheap compared to other articulatory techniques. We developed a smartphone application that can record silent lip video and synthesize speech [31].

The third technique that we investigated within this WP is real-time MRI (rtMRI) of the vocal tract [32]. MRI has not been used before for articulatory-to-acoustic (forward) mapping; although its advantage is that it has a high 'relative' spatial resolution. We trained various DNNs for articulatory-to-speech conversion, using rtMRI as input, in a speaker-specific way (Fig. 4 shows a CNN-LSTM network for this purpose). We used American English speakers of the USC-TIMIT articulatory database [33]. We evaluated the results with objective (Normalized MSE and MCD) and subjective measures (perceptual test) and showed that CNN-LSTM networks are preferred (similarly to the earlier cases of UTI and LIP). Next, neural vocoders were applied for MRI-to-speech reconstruction, and we showed that the approach can successfully reconstruct the gross spectral shape, but more improvements are needed to reproduce the fine spectral details [34].

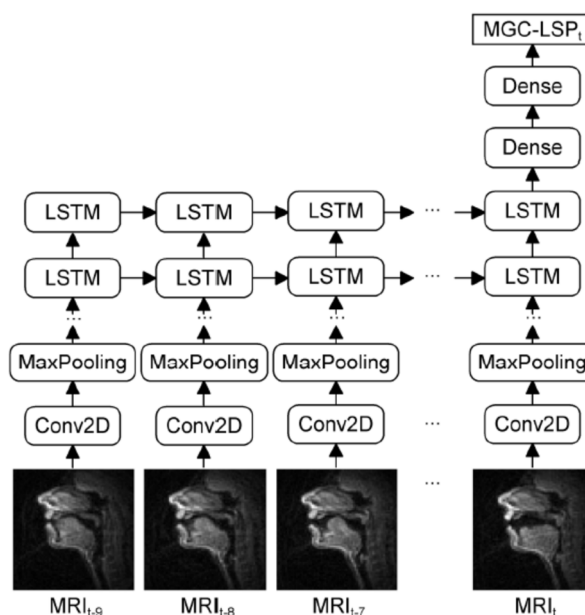


Fig. 4: CNN-LSTM network for MRI-to-speech forward mapping. From [32].

Another interesting finding of our MRI experiments was that 74% of the recordings of speaker `m1' in USC-TIMIT are out of sync [32]. We contacted the developers of the USC-TIMIT database (Asterios Toutios & Shrikanth Narayanan from the University of Southern California, USA) and discussed with them that further investigations are necessary to check this audio-visual synchrony problem. They are open for future collaboration in this field, which can be the basis of a next research grant.

In Sep 2019, as part of the PI's current FK-17 and the other PD-18 projects, we invited Maida Percival (PhD student at the University of Toronto, Canada; expert in ultrasound tongue recordings for linguistic purposes) to Budapest; and this co-operation resulted in numerous presentations [35]–[37]. In 2020, during such ultrasound recording sessions, the PI and colleagues observed a serious methodological issue: a limitation of ultrasound tongue imaging is the transducer misalignment during longer data recording sessions. We presented this problem at various conferences: Interspeech 2020 [38], UltraFest IX [36], ISSP 2020 [39]; and a large number of researchers confirmed that they have similar issues but do not have the solution yet (e.g. Alan Wrench & Pertti Palo, QMU, UK; Kevin Roon & Wei-Rong Chen & Douglas Whalen, Haskins & YALE & CUNY, USA; Michael Pucher, ÖAW, Austria; Judith Dineley, Augsburg University, Germany; Matthew Faytak, UCLA, USA; Sherman Charles & Steven Lulich, Indiana University, USA; Martijn Wieling, University of Groningen, The Netherlands; Kiyoshi Honda, Tianjin University, China – to mention a few). Therefore, further investigations of speaker dependency and solving the above problem of ultrasound transducer misalignment will be definitely useful for the whole speech community, dealing with articulation or speech production research.

Lastly, we conducted cross-speaker experiments, using x-vectors [40]. Our first attempt to apply them in a multi-speaker silent speech framework brought about a marginal reduction in the error rate of the spectral estimation step. Besides, we created an initial feasibility study for text-to-ultrasound prediction [41]. We extended a traditional (vocoder-based) DNN-TTS framework with predicting PCA-compressed ultrasound images, of which the continuous tongue motion can be reconstructed in synchrony with synthesized speech. Articulatory movement prediction from text input can be useful for audiovisual speech synthesis. A specific application is computer-assisted pronunciation training / computer-aided language learning, which can be beneficial for learners of second languages.

Overall, in this WP, we conducted numerous experiments from various aspects (input representations, machine learning approaches, target representations, and session / speaker dependency), which all can help to create practical Silent Speech Interface systems.

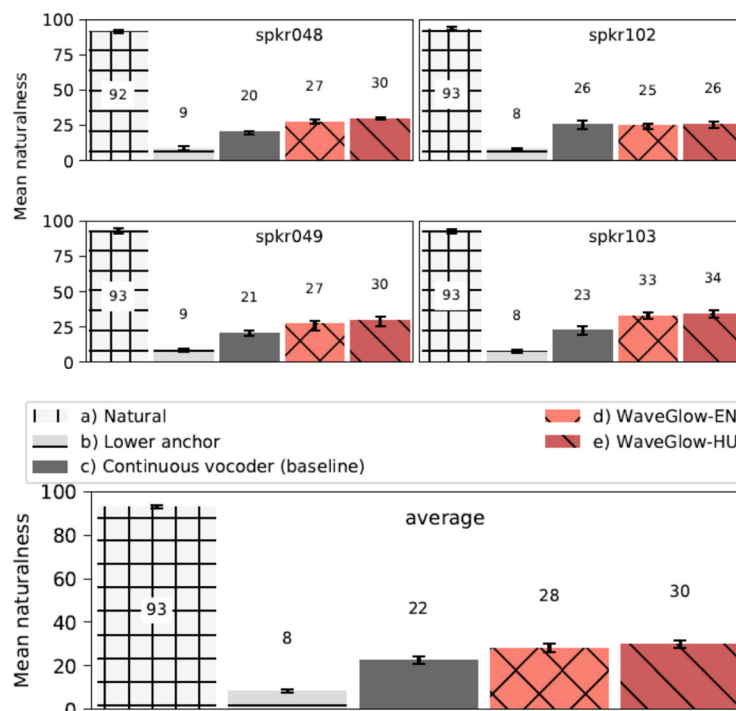


Fig. 5: Results of the subjective evaluation with respect to naturalness, speaker by speaker (top) and average (bottom). The errorbars show the 95% confidence intervals. From [24].

2.4 WP4: Evaluation and testing with users

Within all of the previous tasks, we performed objective evaluations, and listening tests as subjective evaluations. An example for the results of a MUSHRA-like listening test is presented in Fig. 5, which shows that with the WaveGlow neural vocoder higher quality synthesized speech can be achieved than with previous vocoders [24].

We developed prototype systems for Silent Speech Interfaces. Our developed contributions and prototype systems are available open-source:

- <https://github.com/BME-SmartLab/UTI-to-STFT>
- <https://github.com/BME-SmartLab/txt-ult2wav>
- <https://github.com/BME-SmartLab/txt2ult>
- <https://github.com/BME-SmartLab/UTI-to-STFT-Tacotron2>
- <https://github.com/BME-SmartLab/UTI-misalignment>
- <https://github.com/BME-SmartLab/speech2uti>
- <https://github.com/BME-SmartLab/UTI-raw-vs-wedge>
- <https://github.com/malradhi/conTTS>
- <https://github.com/BME-SmartLab/mri2speech>
- <https://github.com/BME-SmartLab/speech2mri>

3 Summary and conclusions

Articulatory-to-acoustic mapping is a novel research field within speech technology, having Silent Speech Interfaces as a long-term potential application. Within this project, we 1) proposed novel methods for recognition-and-synthesis and direct synthesis in the field of SSI, 2) analyzed and compared several articulatory tracking methods (ultrasound tongue imaging, lip video, and vocal tract MRI), 3) contributed to acoustic-to-articulatory inversion. We conducted numerous experiments, including ultrasound-to-speech, ultrasound-to-text, ultrasound-to-F0, lip-to-speech, MRI-to-speech, employed various deep learning architectures (fully connected, 2D and 3D convolutional, recurrent neural networks, multi-task learning, autoencoders, generative adversarial networks), compared continuous and neural vocoders, investigated voice activity detection from ultrasound, and tested ultrasound data representations (raw scanlines / wedge orientation / PCA compression), recorded new Hungarian data and also applied English datasets. We first trained speaker-dependent neural networks, and later we proposed solutions for cross-session and cross-speaker articulation-to-speech synthesis, proceeding towards a practical prototype.

During the four years of the project, the main focus was on the ultrasound modality, but we also investigated video of the lip movement and MRI of the vocal tract. Numerous BSc/MSc/PhD students were involved in the projects, as part of their project laboratory, thesis writing, internship or individual research project. Table 1 summarizes that there were nine related BSc topics, 13 MSc theses, and 6 PhD research topics.

We presented our results at high-ranked conferences (e.g. Interspeech, ISSP, SSW) and published in top journals (e.g. Multimedia Tools and Applications, Computer Speech and Language). Already during the four years of the project, we received 60+ citations. Those appearing in key journals are listed here:

- Csapó & Xu, 2020 [38] is cited by Ribeiro et al. 2021 [42] (Speech Communication)
- Al-Radhi et al., 2020 [10] is cited by Aichinger & Pernkopf 2021 [43] (IEEE/ACM Transactions on Audio, Speech and Language Processing)
- Csapó et al., 2020 [24] is cited by Zhang et al., 2021 [44] (IEEE Access)
- Gosztolya et al., 2020 [23] is cited by Gonzalez-Lopez et al., 2020 [45] (IEEE Access)
- Csapó et al., 2019 [12] is cited by Lee et al., 2020 [46] (Sensors) and Gonzalez-Lopez et al., 2020 [45] (IEEE Access)
- Porras et al., 2019 [5] is cited by Eshky et al., 2021 [47] (Speech Communication) and Shahrehabaki et al., 2021 [48] (IEEE/ACM Transactions on Audio, Speech and Language Processing)
- Gosztolya et al., 2019 [21] is cited by Wang et al., 2021 [49] (Journal of the Acoustical Society of America – Express Letters), Lee et al., 2020 [46] (Sensors) and Gonzalez-Lopez et al., 2020 [45] (IEEE Access)

- Grósz et al., 2018 [19] is cited by Zhang et al., 2021 [44] (IEEE Access), Wang et al., 2021 [49] (Journal of the Acoustical Society of America – Express Letters) and Parlak&Altun, 2021 [50] (Mathematical Problems in Engineering)
- Tóth et al., 2018 [4] is cited by Zhang et al., 2021 [35] (IEEE Access), Gonzalez-Lopez et al., 2020 [36] (IEEE Access), and Wu&Weng, 2021 [51] (Neural Networks)

Table 1: Students involved in the FK-17 project.

Name	Type / Year	Topic
Rémi Balandras	BSc project / 2017	acoustic-to-articulatory inversion (AAI) using modern machine learning (French & ultrasound)
Léa Desse	BSc project / 2017	AAI using modern machine learning (French & ultrasound)
Chrysogone Paolo	BSc project / 2017	AAI using modern machine learning (French & ultrasound)
Eloi Moliner	BSc thesis / 2017-2018	ultrasound-to-speech using CNNs & LSTMs
Akif Alic	MSc thesis / 2017-2018	voice conversion using machine learning
Balaton Tamás	BSc project / 2018	Variational AutoEncoders & ultrasound
Blaskó Gergő	BSc thesis / 2018	AAI based on ultrasound tongue imaging, using deep learning
Dagoberto Porras Plata	MSc internship / 2018-2019	ultrasound & acoustic-to-articulatory inversion
Makrai Márton	PhD individual research topic / 2018-2019	hyperparameter optimization of deep neural networks for ultrasound-to-speech synthesis
Nadia Hajjej	MSc thesis / 2018-2019	application of Generative Adversarial Networks (GANs) for processing of ultrasound images
Amarsaikhan Nasantogtok	MSc thesis / 2018-2019	applying CycleGANs for ultrasound-speech conversion
Khorkova Mariia	MSc thesis / 2018-2020	lip-to-speech synthesis using DNNs
Ráczi Bianka	BSc thesis / 2018-2019	lip-to-speech synthesis using CNNs and RNNs
Varga Kristóf	BSc thesis / 2018-2019	lip-to-speech synthesis using face tracking
Bárány Bálint	MSc thesis / 2018-2020	ultrasound-to-speech synth., iterative data loading
Arthur Viktor	MSc thesis / 2018-2022	vid-to-speech mobile application
Maida Percival	PhD individual research topic / 2019-2020	ultrasound session & speaker dependency, manual contour tracking, transducer auto-rotation, linguistic aspects
Amin Shandiz	PhD individual research topic / 2019-2021	data augmentation for UTI; speaker dependency, voice activity detection from articulatory data
José Lopez	PhD individual research topic / 2019-2021	testing neural network architectures for ultrasound-to-speech, hyperparameter optimization
Mohammed Al-Radhi	PhD individual research topic / 2018-2022	spectral filtering for text-to-speech, voice conversion, enhanced spectral filtering with articulatory data
Törner Márton	MSc project / 2020	DNN types for ultrasound-to-speech
Pengyu Dai	MSc thesis / 2019-2020	ultrasound-to-F0
Dóka Zsolt	BSc thesis / 2019-2020	lip-to-speech with MagPhase vocoder
Szokoly Virág	MSc project / 2021	DNN types for ultrasound-to-speech and TTS
Yide Yu	MSc thesis / 2020-2021	MRI-to-speech using neural vocoders
Wu Liang	MSc thesis / 2020-2021	acoustic-to-articulatory inversion using ultrasound
Rysbekova Aliia	MSc thesis / 2020-2021	Comparison of UTI-EMA-MRI-LIP video and its application for DNN-based articulation-to-speech
Mohammad Areej Mohammad Mousa	MSc thesis / 2020-2022	Ultrasound Tongue Imaging for Silent Speech Interfaces using deep learning
Ali Raheem Mandeel	PhD individual research topic / 2020-2022	speaker adaptation for text-to-speech synthesis

References

- [1] B. Denby, T. Schultz, K. Honda, T. Hueber, J. M. Gilbert, and J. S. Brumberg, “Silent speech interfaces,” *Speech Commun.*, vol. 52, no. 4, pp. 270–287, 2010.
- [2] M. Wand, J. Koutník, and J. Schmidhuber, “Lipreading with long short-term memory,” in *Proc. ICASSP*, 2016, pp. 6115–6119.
- [3] B. Cao, M. Kim, J. R. Wang, J. Van Santen, T. Mau, and J. Wang, “Articulation-to-Speech Synthesis Using Articulatory Flesh Point Sensors’ Orientation Information,” in *Proc. Interspeech*, 2018, pp. 3152–3156.
- [4] L. Tóth, G. Gosztolya, T. Grósz, A. Markó, and T. G. Csapó, “Multi-Task Learning of Phonetic Labels and Speech Synthesis Parameters for Ultrasound-Based Silent Speech Interfaces,” in *Proc. Interspeech*, 2018, pp. 3172–3176.
- [5] D. Porras, A. Sepulveda-Sepulveda, and T. G. Csapo, “DNN-based Acoustic-to-Articulatory Inversion using Ultrasound Tongue Imaging,” in *Proceedings of the International Joint Conference on Neural Networks*, 2019, vol. 2019-July.
- [6] T. G. Csapó and A. Sepúlveda, “Ultrasound Tongue Image Generation for Acoustic-to-Articulatory Inversion using Convolutional and Recurrent Deep Neural Networks,” *Submitt. to Multimed. Tools Appl.*, 2021.
- [7] T. G. Csapó, “Speaker dependent acoustic-to-articulatory inversion using real-time MRI of the vocal tract,” in *Proc. Interspeech*, 2020, pp. 3720–3724.
- [8] M. S. Al-Radhi, T. G. Csapó, and G. Németh, “Continuous vocoder applied in deep neural network based voice conversion,” *Multimed. Tools Appl.*, vol. 78, no. 23, 2019.
- [9] M. S. Al-Radhi, O. Abdo, T. G. Csapó, S. Abdou, G. Németh, and M. Fashal, “A continuous vocoder for statistical parametric speech synthesis and its evaluation using an audio-visual phonetically annotated Arabic corpus,” *Comput. Speech Lang.*, vol. 60, p. 101025, Mar. 2020.
- [10] M. S. Al-Radhi, T. G. Csapó, and G. Németh, “Continuous Noise Masking Based Vocoder for Statistical Parametric Speech Synthesis,” *IEICE Trans. Inf. Syst.*, vol. E103.D, no. 5, pp. 1099–1107, May 2020.
- [11] M. S. Al-Radhi, T. G. Csapo, and G. Nemeth, “Parallel voice conversion based on a continuous sinusoidal model,” in *2019 10th International Conference on Speech Technology and Human-Computer Dialogue, SpeD 2019*, 2019.
- [12] T. G. Csapó *et al.*, “Ultrasound-based Silent Speech Interface Built on a Continuous Vocoder,” in *Proc. Interspeech*, 2019, pp. 894–898.
- [13] M. S. Al-Radhi, T. G. Csapó, and G. Németh, “Noise and acoustic modeling with waveform generator in text-to-speech and neutral speech conversion,” *Multimed. Tools Appl.*, vol. 80, no. 2, pp. 1969–1994, Jan. 2021.
- [14] M. S. Al-Radhi, T. G. Csapó, C. Zainkó, and G. Németh, “Continuous wavelet vocoder-based decomposition of parametric speech waveform synthesis,” in *Proc. Interspeech*, 2021, pp. 3266–3270.
- [15] M. S. Al-Radhi, T. G. Csapó, and G. Németh, “conTTS: Text-to-Speech Application using a Continuous Vocoder,” in *Proc. ISSP*, 2020.
- [16] T. G. Csapó, T. Grósz, G. Gosztolya, L. Tóth, and A. Markó, “DNN-Based Ultrasound-to-Speech Conversion for a Silent Speech Interface,” in *Proc. Interspeech*, 2017, pp. 3672–3676.
- [17] E. Moliner and T. G. Csapó, “Ultrasound-based silent speech interface using convolutional and recurrent neural networks,” *Acta Acust. united with Acust.*, vol. 105, no. 4, pp. 587–590, 2019.
- [18] L. Tóth and A. H. Shandiz, “3D Convolutional Neural Networks for Ultrasound-Based Silent Speech Interfaces,” in *Proc. ICAISC*, 2020.
- [19] T. Grósz, G. Gosztolya, L. Tóth, T. G. Csapó, and A. Markó, “F0 Estimation for DNN-Based Ultrasound Silent Speech Interfaces,” in *Proc. ICASSP*, 2018, pp. 291–295.
- [20] P. Dai, M. S. Al-Radhi, and T. G. Csapo, “Effects of F0 Estimation Algorithms on Ultrasound-Based Silent Speech Interfaces,” in *SpeD*, 2021, pp. 47–51.
- [21] G. Gosztolya, Á. Pintér, L. Tóth, T. Grósz, A. Markó, and T. G. Csapó, “Autoencoder-Based Articulatory-to-Acoustic Mapping for Ultrasound Silent Speech Interfaces,” in *International Joint Conference on Neural Networks*, 2019.
- [22] A. Honarmandi Shandiz and L. Tóth, “Voice Activity Detection for Ultrasound-Based Silent Speech Interfaces Using Convolutional Neural Networks,” *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 12848 LNAI, pp. 499–510, Sep. 2021.
- [23] G. Gosztolya, T. Grósz, L. Tóth, A. Markó, and T. G. Csapó, “Applying DNN Adaptation to Reduce the Session Dependency of Ultrasound Tongue Imaging-Based Silent Speech Interfaces,” *Acta Polytech. Hungarica*, vol. 17, no. 7, pp. 109–124, 2020.
- [24] T. G. Csapó, C. Zainkó, L. Tóth, G. Gosztolya, and A. Markó, “Ultrasound-based Articulatory-to-Acoustic Mapping with WaveGlow Speech Synthesis,” in *Proc. Interspeech*, 2020, pp. 2727–2731.

- [25] N. Hajjaj and T. G. Csapó, “Realistic Ultrasound Tongue Image Synthesis using Generative Adversarial Networks,” *Beszédtudomány - Speech Sci.*, vol. 1, no. 1, pp. 7–21, 2020.
- [26] T. G. Csapó, G. Gosztolya, L. Tóth, A. H. Shandiz, and A. Markó, “Optimizing the Ultrasound Tongue Image Representation for Residual Network-based Articulatory-to-Acoustic Mapping,” *Submitt. to Multimed. Tools Appl.*, 2021.
- [27] C. Zainkó *et al.*, “Adaptation of Tacotron2-based Text-To-Speech for Articulatory-to-Acoustic Mapping using Ultrasound Tongue Imaging,” in *Proc. ISCA SSW11*, 2021, pp. 54–59.
- [28] T. G. Csapó, L. Tóth, G. Gosztolya, and A. Markó, “Speech Synthesis from Text and Ultrasound Tongue Image-based Articulatory Input,” in *Proc. ISCA SSW11*, 2021, pp. 31–36.
- [29] B. Rácz and T. G. Csapó, “Ajakvideó alapú beszédszintézis konvolúciós és rekurrens mély neurális hálózatokkal,” *Beszédtudomány -- Speech Sci.*, pp. 57–72, 2020.
- [30] M. Cooke, J. Barker, S. Cunningham, and X. Shao, “An audio-visual corpus for speech perception and automatic speech recognition,” *J. Acoust. Soc. Am.*, vol. 120, no. 5, pp. 2421–2424, Nov. 2006.
- [31] F. V. Arthur and T. G. Csapó, “Towards a practical lip-to-speech conversion system using deep neural networks and mobile application frontend,” in *2nd International Conference on Artificial Intelligence and Computer Vision (AICV2021)*, 2021.
- [32] T. G. Csapó, “Speaker dependent articulatory-to-acoustic mapping using real-time MRI of the vocal tract,” in *Proc. Interspeech*, 2020, pp. 2722–2726.
- [33] S. Narayanan *et al.*, “Real-time magnetic resonance imaging and electromagnetic articulography database for speech production research (TC),” *J. Acoust. Soc. Am.*, vol. 136, no. 3, pp. 1307–1311, Sep. 2014.
- [34] Y. Yu, A. Honarmandi Shandiz, and L. Tóth, “Reconstructing Speech from Real-Time Articulatory MRI Using Neural Vocoders,” in *Proc. EUSIPCO*, 2021, pp. 945–949.
- [35] M. Percival, T. G. Csapó, M. Bartók, A. Deme, T. E. Grácsi, and A. Markó, “Gemination as fortition? Articulatory data from Hungarian,” in *LabPhon*, 2020.
- [36] M. Percival, T. G. Csapó, M. Bartók, A. Deme, T. E. Grácsi, and A. Markó, “Ultrasound imaging of Hungarian geminates,” in *UltraFest IX*, 2020.
- [37] M. Percival, T. G. Csapó, M. Bartók, A. Deme, T. E. Grácsi, and A. Markó, “Tongue root and voicing in Hungarian singleton and geminate obstruents,” in *12th International Seminar on Speech Production*, 2020.
- [38] T. G. Csapó and K. Xu, “Quantification of Transducer Misalignment in Ultrasound Tongue Imaging,” in *Proc. Interspeech*, 2020, pp. 3735–3739.
- [39] T. G. Csapó, K. Xu, A. Deme, T. E. Grácsi, and A. Markó, “Transducer Misalignment in Ultrasound Tongue Imaging,” in *12th International Seminar on Speech Production*, 2020.
- [40] A. H. Shandiz, L. Toth, Gabor Gosztolya, A. Markó, and T. Gabor Csapó, “Neural speaker embeddings for ultrasound-based silent speech interfaces,” in *Proc. Interspeech*, 2021, vol. 1, pp. 151–155.
- [41] T. G. Csapó, “Extending Text-to-Speech Synthesis with Articulatory Movement Prediction using Ultrasound Tongue Imaging,” in *Proc. ISCA SSW11*, 2021, pp. 7–12.
- [42] M. S. Ribeiro, J. Cleland, A. Eshky, K. Richmond, and S. Renals, “Exploiting ultrasound tongue imaging for the automatic detection of speech articulation errors,” *Speech Commun.*, vol. 128, pp. 24–34, Apr. 2021.
- [43] P. Aichinger and F. Pernkopf, “Synthesis and Analysis-By-Synthesis of Modulated Diplophonic Glottal Area Waveforms,” *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 914–926, 2021.
- [44] J. Zhang, P. Roussel, and B. Denby, “Creating Song from Lip and Tongue Videos with a Convolutional Vocoder,” *IEEE Access*, vol. 9, pp. 13076–13082, 2021.
- [45] J. A. Gonzalez-Lopez, A. Gomez-Alanis, J. M. Martin Donas, J. L. Perez-Cordoba, and A. M. Gomez, “Silent Speech Interfaces for Speech Restoration: A Review,” *IEEE Access*, vol. 8, pp. 177995–178021, Sep. 2020.
- [46] W. Lee, J. J. Seong, B. Ozlu, B. S. Shim, A. Marakhimov, and S. Lee, “Biosignal sensors and deep learning-based speech recognition: A review,” *Sensors*, vol. 21, no. 4, pp. 1–22, 2021.
- [47] A. Eshky, J. Cleland, M. S. Ribeiro, E. Sugden, K. Richmond, and S. Renals, “Automatic audiovisual synchronisation for ultrasound tongue imaging,” *Speech Commun.*, vol. 132, pp. 83–95, Sep. 2021.
- [48] A. Sabzishahrehabaki, G. Salvi, T. K. Svendsen, and S. M. Siniscalchi, “Acoustic-to-Articulatory Mapping with Joint Optimization of Deep Speech Enhancement and Articulatory Inversion Models,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, pp. 1–1, 2021.
- [49] H. Wang, P. Roussel, and B. Denby, “Improving ultrasound-based multimodal speech recognition with predictive features from representation learning,” *JASA Express Lett.*, vol. 1, no. 1, p. 015205, Jan. 2021.
- [50] C. Parlak and Y. Altun, “Harmonic Differences Method for Robust Fundamental Frequency Detection in Wideband and Narrowband Speech Signals,” *Math. Probl. Eng.*, vol. 2021, 2021.
- [51] X. Wu and J. Weng, “Learning to recognize while learning to speak: Self-supervision and developing a speaking motor,” *Neural Networks*, vol. 143, pp. 28–41, Nov. 2021.