

A Multimodal Corpus of Expert Gaze and Behavior during Phonetic Segmentation Tasks

Arif Khan¹⁻³, Ingmar Steiner^{1,2}, Yusuke Sugano⁴, Andreas Bulling^{1,5}, Ross Macdonald⁶

¹Multimodal Computing and Interaction, Saarland University

³Saarbrücken Graduate School of Computer Science

⁵Max Planck Institute for Informatics, Saarbrücken, Germany

²DFKI Language Technology Lab, Saarbrücken

⁴Osaka University, Japan

⁶University of Manchester, UK

1 Introduction

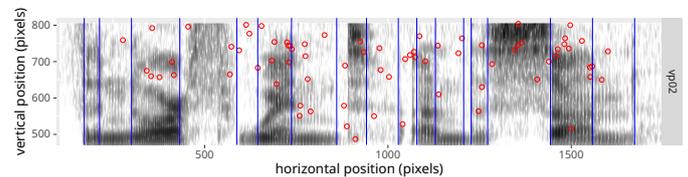
- Phonetic segmentation is the process of splitting speech into distinct phonetic units.
- Automatic phonetic segmentation should replicate the precision of human segmentation as closely as

- possible, but segmentation behavior data is scarce.
- This corpus captures human segmentation behavior by recording phonetician's gaze with an eyetracker, along with other relevant modalities.

2 Data recording

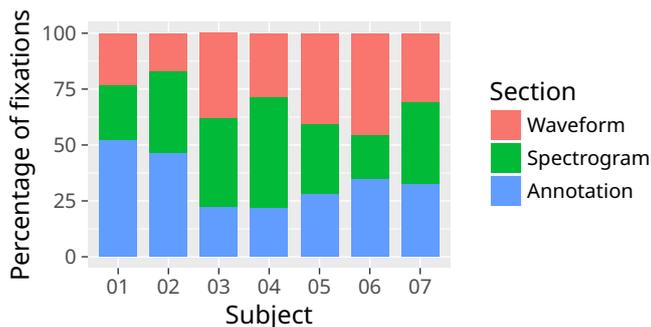
- Segmentation of a 46 s audio recording (48 kHz sampling rate).
- The audio was segmented using the *Praat* software.
- Segmenter gaze was recorded with a Tobii TX300 eyetracker.
- Playback audio, webcam video, and screen contents were also recorded.
- The *Praat* UI state and final manual segmentation were also stored.

3 Gaze on spectrogram

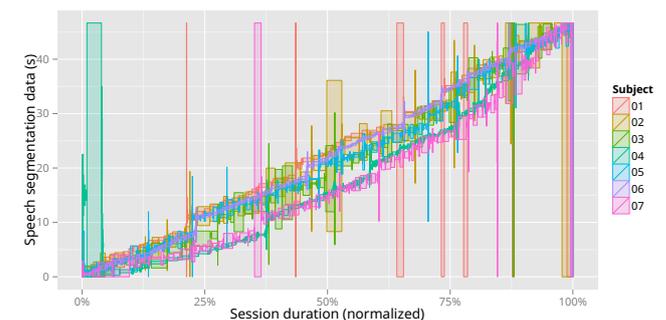


- Spectrogram of the utterance "The North Wind and the Sun".
- Fixations are rendered as red circles.
- Blue lines represent the manual boundaries placed by the participant.

4 Fixation regions



5 Segmentation progress



6 Conclusion

- Recorded a multimodal corpus of segmentation behavior data from phonetic experts.
- All relevant information sources were recorded, e.g., gaze, playback audio, video, and screen recording.
- The processed data is released under a Creative Commons license and publicly available on GitHub: <https://git.io/eyeseg-data>.

7 Acknowledgements

- We are grateful to our participants for their time and valuable feedback.
- This project was funded by the German Research Foundation (DFG) under grant EXC 284.

