

1. Background

- Automatic Speech Recognition has its limitations.
- The TDNN-BLSTM architecture improved the Phoneme Recognition for Dutch speech.¹
- This architecture has not been tested on English speech.

/hat/ vs. /kat/
/seik/ vs. /'sa:ki/

2. Research questions

How does the TDNN-BLSTM architecture perform on English **read** and **spontaneous** speech?

Comparing results **quantitatively**:

- What is the Phoneme Error Rate for TDNN-BLSTM on spontaneous and on read speech?

Comparing results **qualitatively**:

- What phonemes have a large PER difference between read and spontaneous speech?

3. Methodology

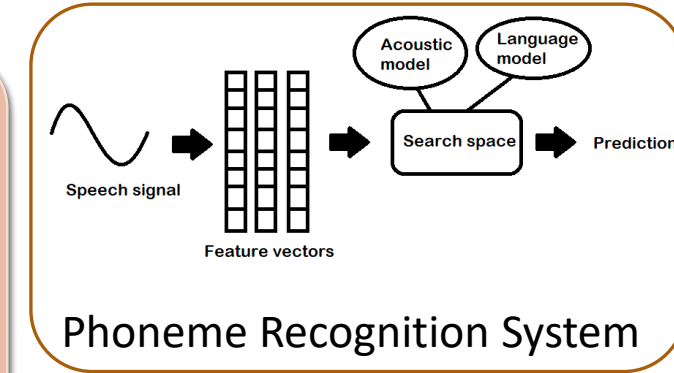
- Using the Timit and Buckeye corpora
- Preparing the data for Kaldi framework
- Training and testing the TDNN-BLSTM model
- Focus on layer dimension, epochs and learning rate
- Evaluating results with PER metric
- Evaluating results based on confusion matrix

4. Results

31.78% PER on read speech
54.03% PER on spontaneous speech

Phonemes in read speech are recognized better than in spontaneous speech

Comparing to research on Dutch PR¹, parallel research² and literature values³:
TDNN-BLSTM does not perform as well as other acoustic models for both spontaneous and read speech.



Some terms explained:

TDNN-BLSTM:

Projected Bidirectional Long Short-Term Memory Time Delayed Neural Network

Read speech: (Timit corpus)

Participants read predetermined sentences

Spontaneous speech: (Buckeye corpus)

Participants have an informal conversation

PER:

Phoneme Error Rate

¹ R. Levenbach, "Phon times: Improving Dutch phoneme recognition," Master's thesis, 2021.

² J. van der Tang, "Evaluation of phoneme recognition through TDNN-OPGRU on Mandarin speech." 2021; G. Genkov, "Training and testing the TDNN-OPGRU acoustic model on English read and spontaneous speech." 2021; M. Chiroșca, "Evaluating the performance of the TDNN-BLSTM on Mandarin read and spontaneous speech." 2021.

³ M. Ravanelli, P. Brakel, M. Omologo, and Y. Bengio, "Light gated recurrent units for speech recognition," IEEE Transactions on Emerging Topics in Computational Intelligence, vol. 2, pp. 92–102, Apr. 2018;

R. Qader, G. Lecorvé, D. Lolive, and P. Sébillot, "Probabilistic Speaker Pronunciation Adaptation for Spontaneous Speech Synthesis Using Linguistic Features," in International Conference on Statistical Language and Speech Processing (SLSP), (Budapest, Hungary), pp. 229–241, Nov. 2015.

¹ R. Levenbach, “Phon times: Improving Dutch phoneme recognition,” Master’s thesis, 2021.

² G. Genkov, “Training and testing the TDNN-OPRGU acoustic model on English read and spontaneous speech.” 2021.

³ M. Chiroșca, “Evaluating the performance of the TDNN-BLSTM on Mandarin read and spontaneous speech.” 2021.

⁴ J. van der Tang, “Evaluation of phoneme recognition through TDNN-OPGRU on Mandarin speech.” 2021.

⁵ M. Ravanelli, P. Brakel, M. Omologo, and Y. Bengio, “Light gated recurrent units for speech recognition,” IEEE Transactions on Emerging Topics in Computational Intelligence, vol. 2, pp. 92–102, Apr. 2018.

⁶ R. Qader, G. Lecorvé, D. Lolive, and P. Sébillot, “Probabilistic Speaker Pronunciation Adaptation for Spontaneous Speech Synthesis Using Linguistic Features,” in International Conference on Statistical Language and Speech Processing (SLSP), (Budapest, Hungary), pp. 229–241, Nov. 2015.