

Objectives

- Compare how accurately machine learning techniques can estimate from acoustic data:
 - vocal tract constrictions and their locations (tract variables)
 - the location of various articulators (pellets)
- Create gestural specifications for a spontaneous speech database

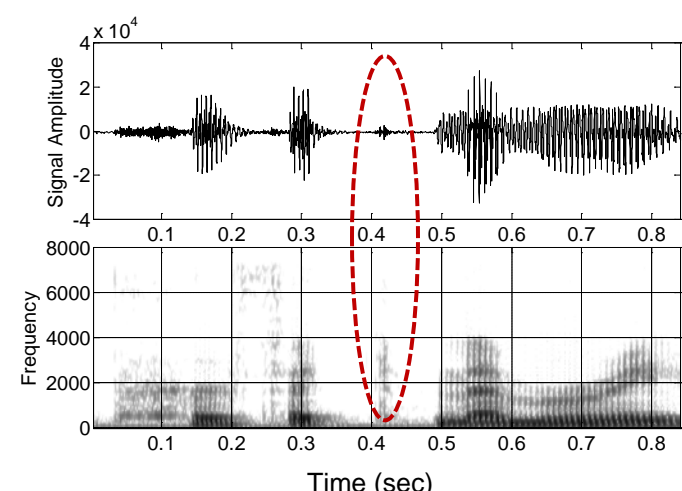
Motivation

Current ASR systems:

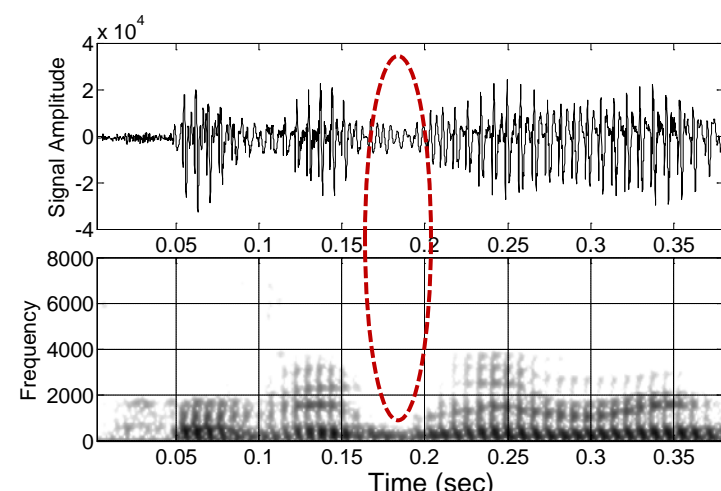
- Based on phones and assume phones to be distinctive regions



- Need to impose limitations in the recognition task in order to handle coarticulation



"perfect-memory" clearly articulated

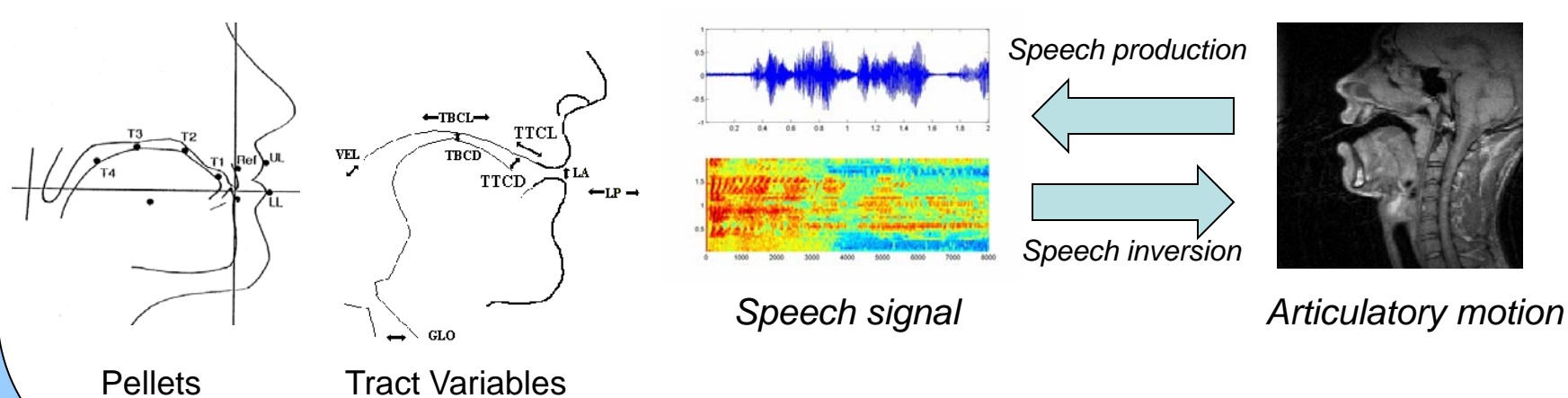


"perfect-memory" quickly articulated

Our Approach

- Use articulatory information to better model coarticulation:

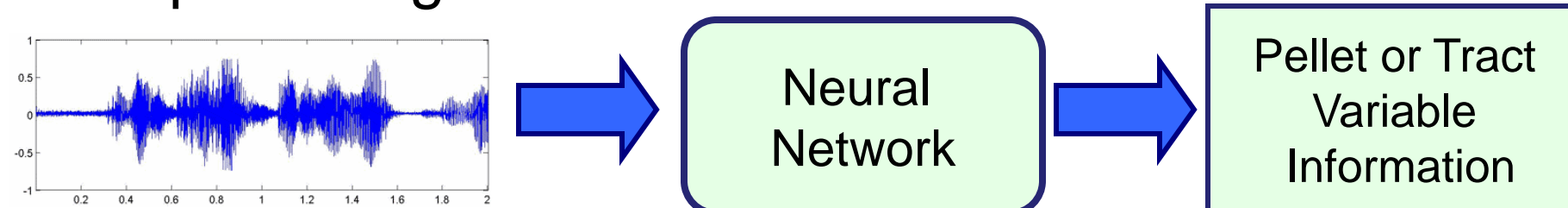
- 7 pellets or 8 tract variables
- Speech Inversion Technique



Methodology

•Task 1: Speech Inversion

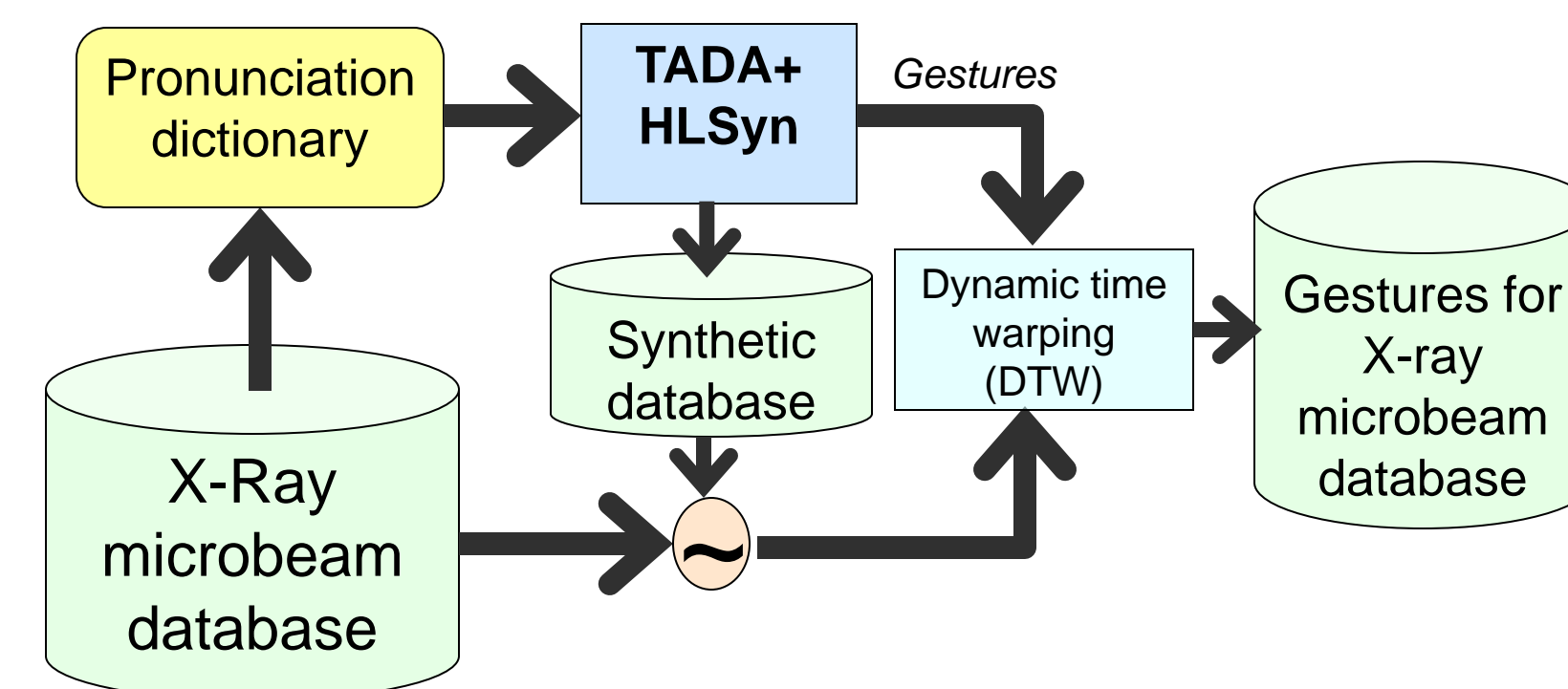
- Train artificial neural networks (ANNs) to estimate tract variables and pellet trajectories given a speech signal



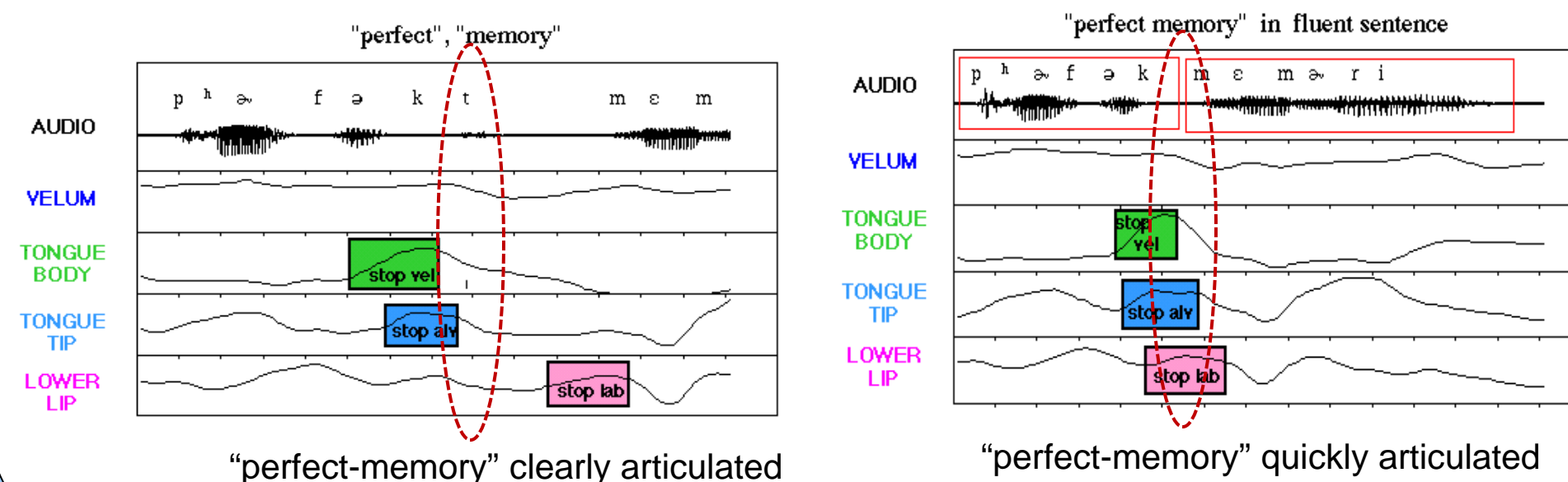
•Task 2: Gestural Modeling

- Gestures are constriction actions along the vocal tract and they are defined by dynamic parameters

•Procedure



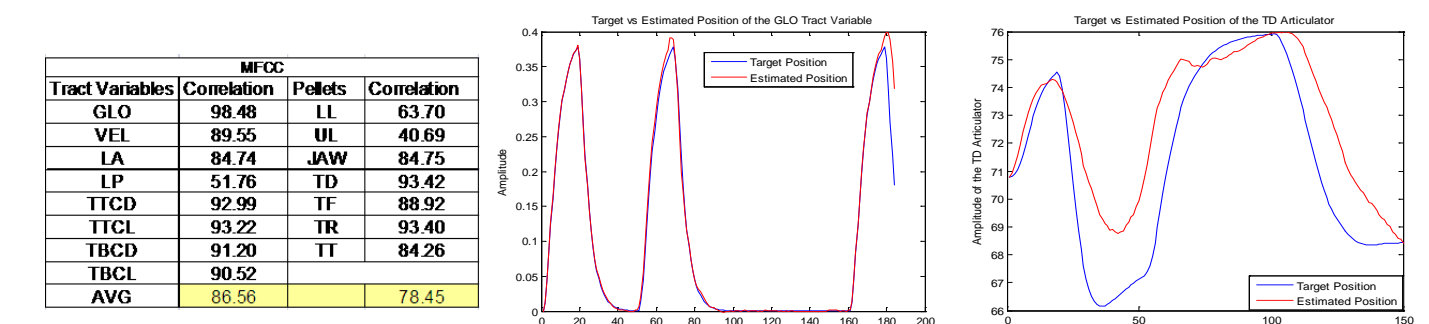
- Gestures are invariant which can account for coarticulation



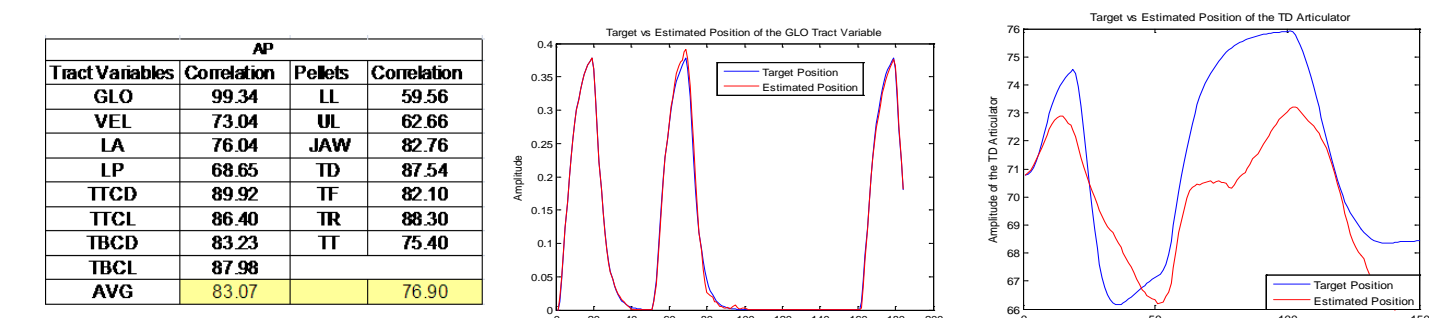
Results

•Neural Network Training Results

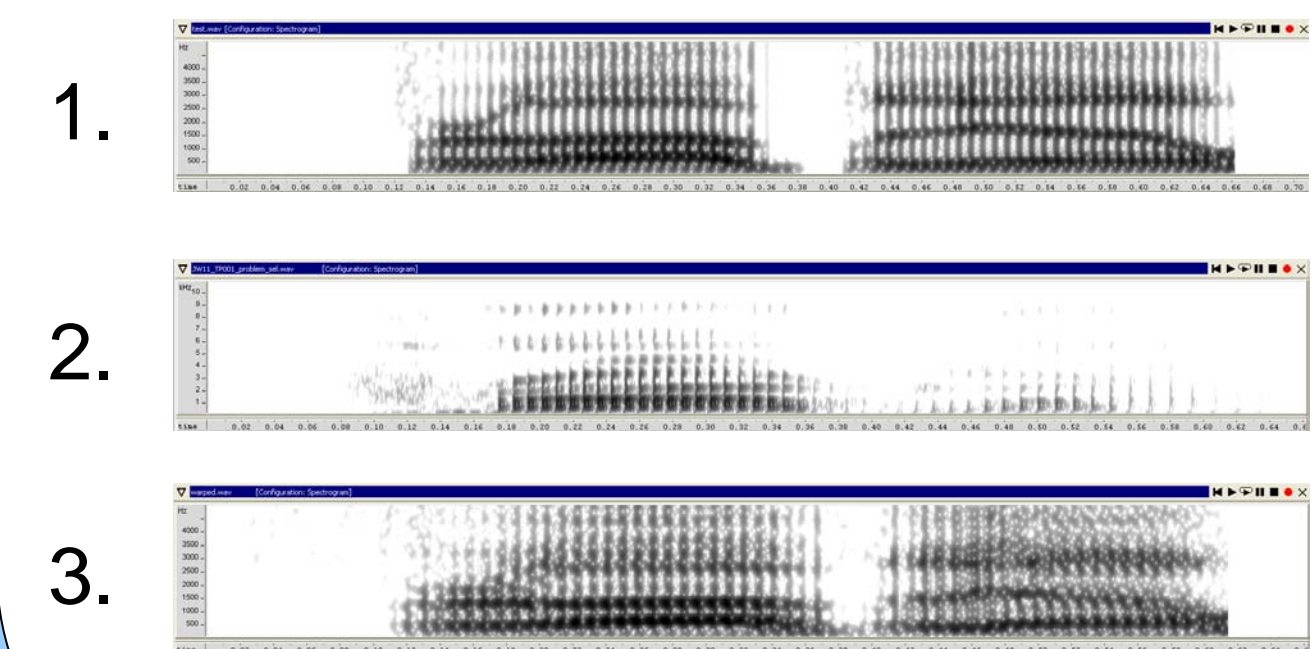
- Using Mel-Frequency Cepstral Coefficients (MFCCs)



- Using Acoustic Parameters (APs)



•Dynamic Time Warping Results



1. Synthetic 2. Natural 3. Warped

Conclusions

- Estimated the tract variables more accurately than the pellets using neural networks
- Warped the synthetic speech signal to the natural speech signal
- Obtained the gestures for the natural speech from the warped synthetic speech

Acknowledgments

- This work was supported by NSF CISE award #0755224