



Aalto University

# UTILIZING GLOTTAL SOURCE PULSE LIBRARY FOR GENERATING IMPROVED EXCITATION SIGNAL FOR HMM-BASED SPEECH SYNTHESIS

Tuomo Raitio<sup>1</sup>, Antti Suni<sup>2</sup>, Hannu Pulakka<sup>1</sup>, Martti Vainio<sup>2</sup>, and Paavo Alku<sup>1</sup>

<sup>1</sup>Department of Signal Processing and Acoustics, Aalto University

<sup>2</sup>Department of Speech Sciences, University of Helsinki

ICASSP 2011

# Contents

- I. Background
- II. Human speech production
- III. Speech synthesis system
- IV. Results and samples

# I. Background

- The ultimate goal of text-to-speech (TTS) is to generate natural sounding expression from arbitrary text
- Two major TTS trends:

## Unit selection

- Based on concatenating prerecorded acoustical units
- Yields (almost) natural quality
- Poor adaptability to speaking styles, speaker characteristics and emotions

## Statistical

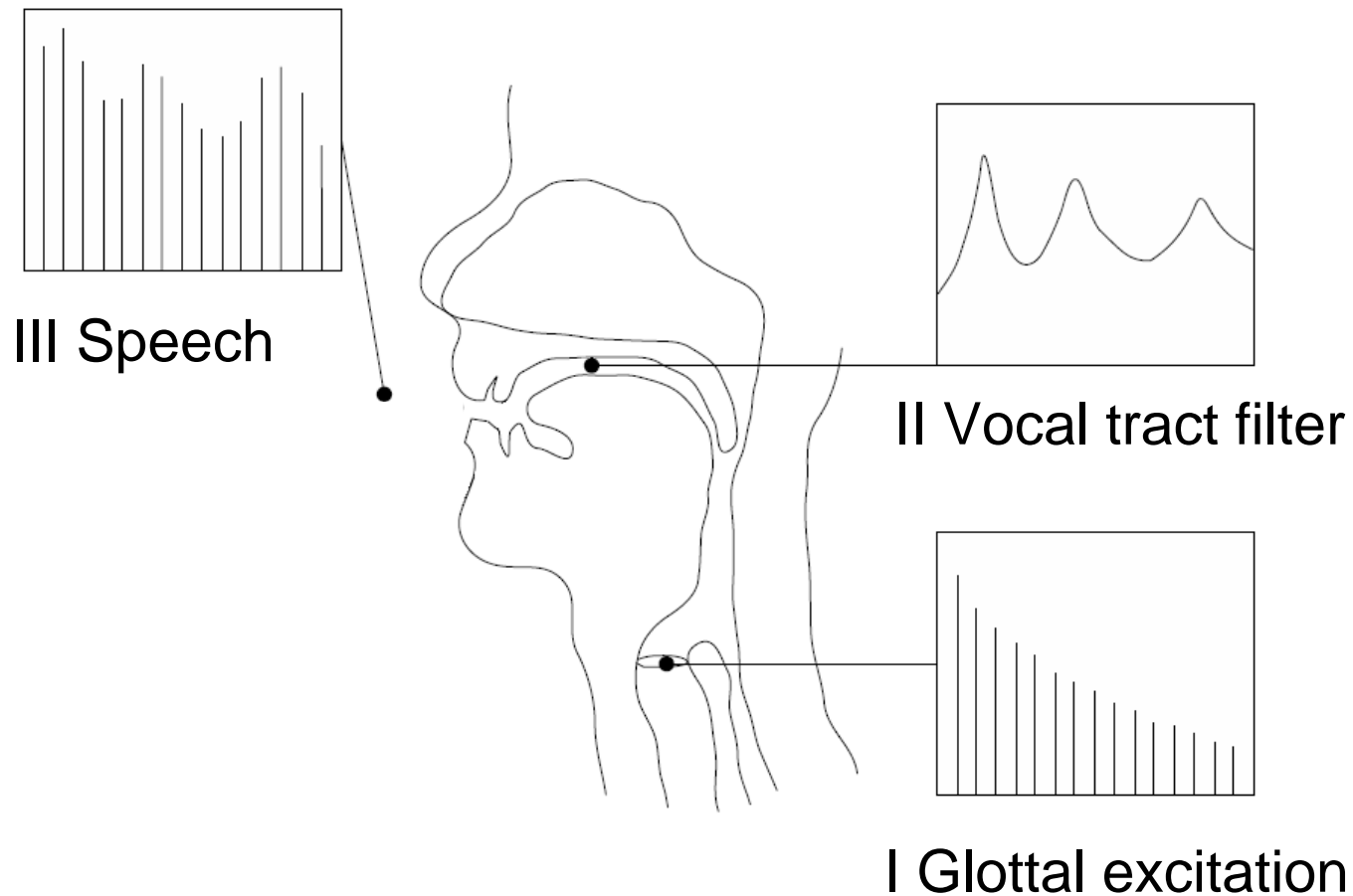
- Based on modeling speech parameters with Hidden Markov Models (HMMs)
- Better adaptability to speaking styles, speaker characteristics and emotions

**Problem:** Current HMM-based synthesizers suffer from degraded naturalness in speech quality

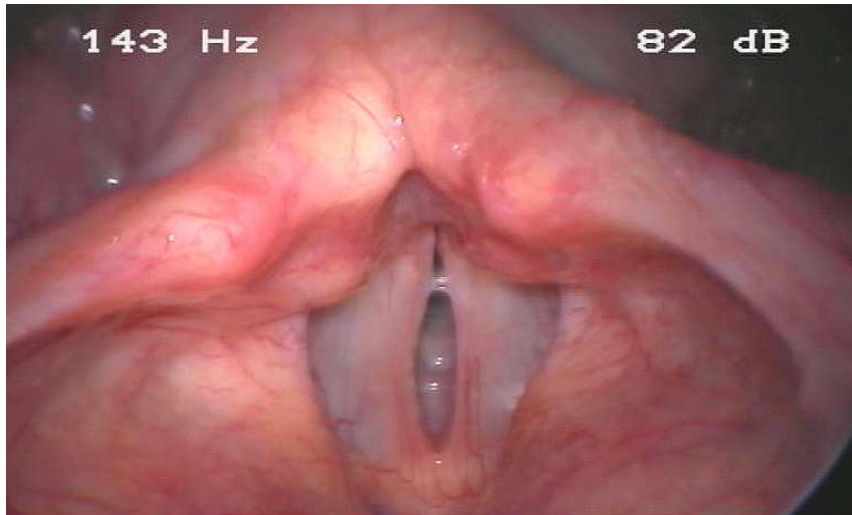
### **Our approach:**

1. Speech is decomposed into the glottal source signal and the vocal tract transfer function
2. Glottal source is further decomposed into several parameters and a glottal pulse library
3. Parameters are modeled in HMMs
4. In synthesis, source signal is reconstructed from the selected glottal pulses and the filtered with the vocal tract filter to create speech

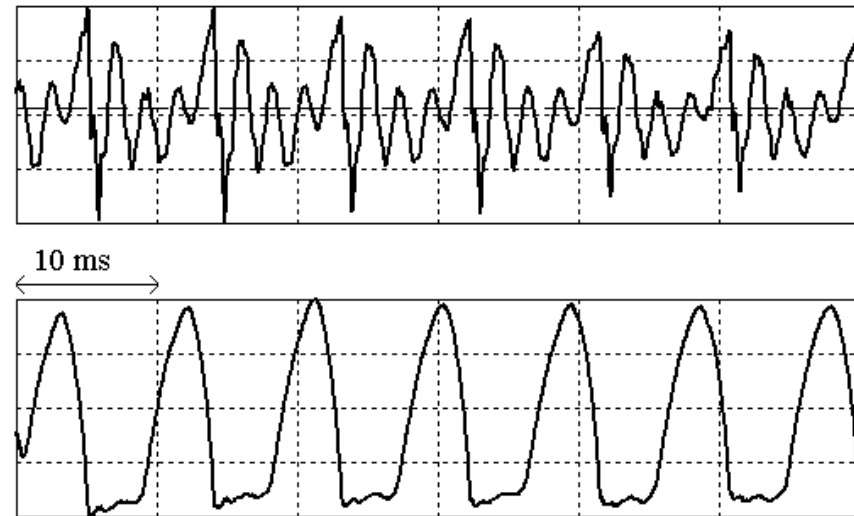
## II. Speech Production Mechanism



# Glottal Source



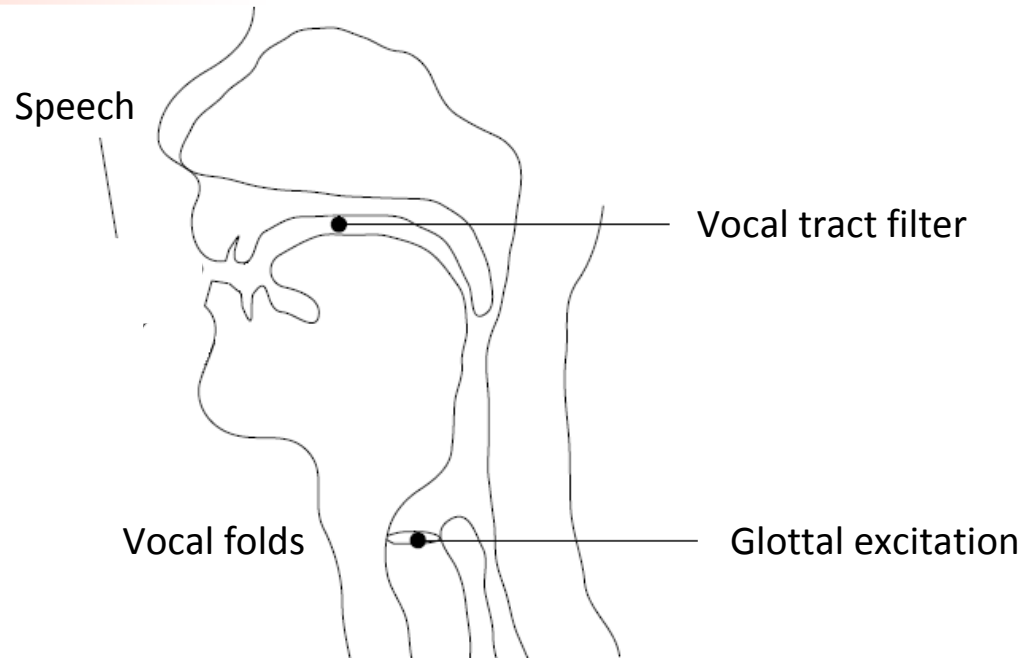
Vibrating vocal folds.



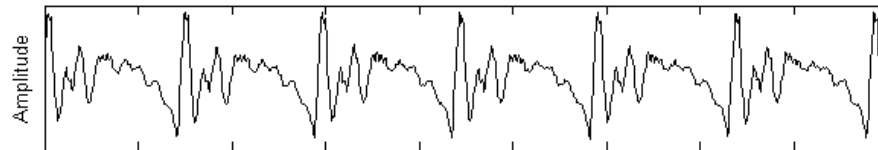
Speech pressure waveform (upper panel) and estimated glottal excitation (lower panel).

# Glottal Inverse Filtering

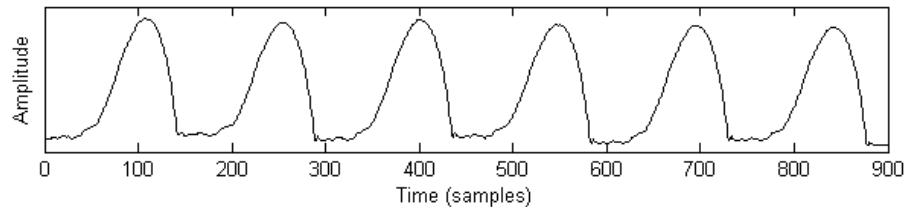
Glottal inverse filtering estimates the glottal flow and the vocal tract filter from a speech signal



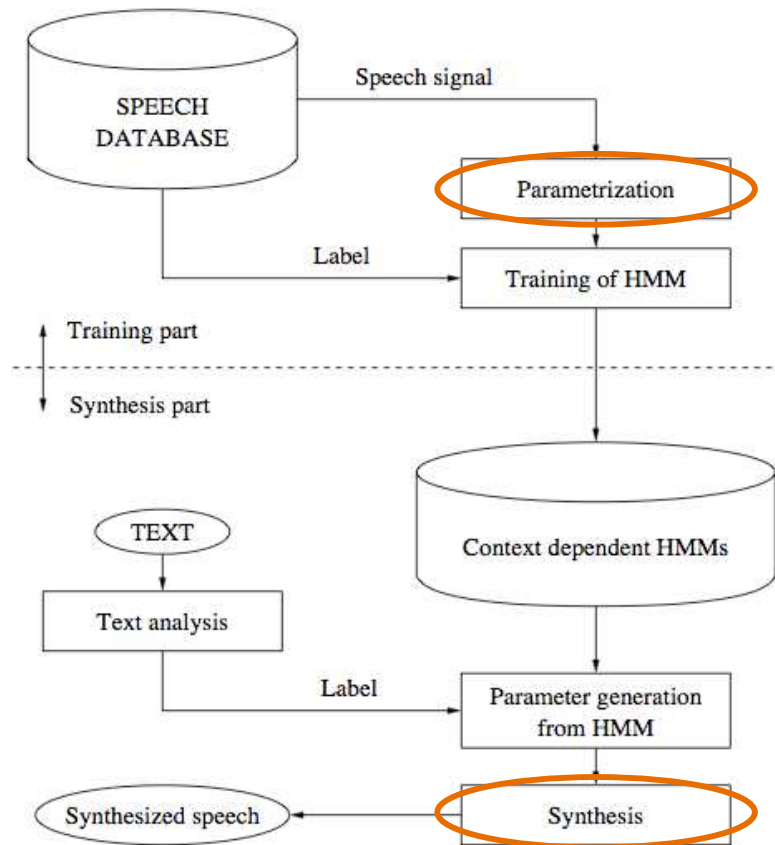
Speech signal



Estimated glottal flow signal

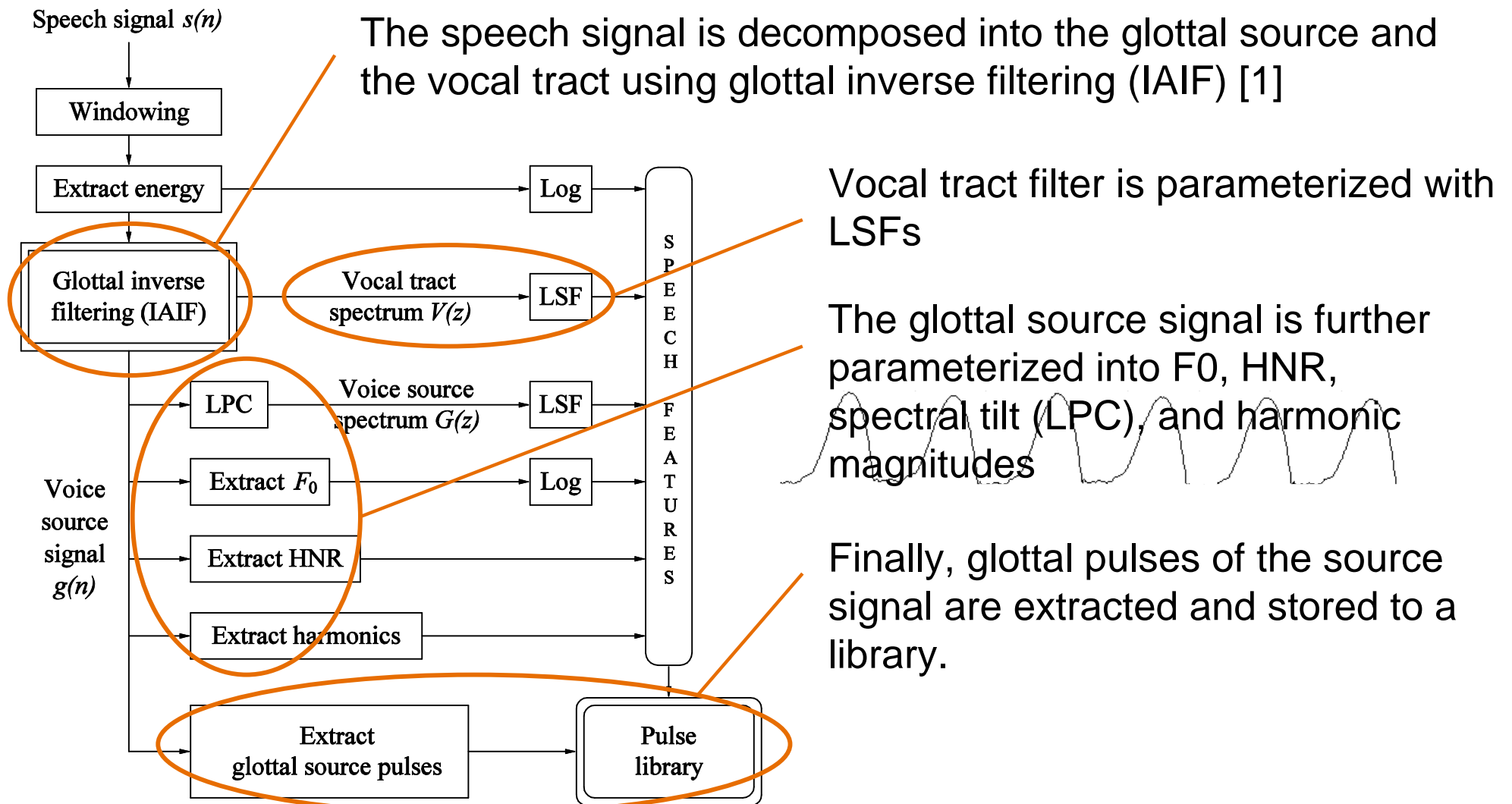


# III. Speech Synthesis System





# Speech Parameterization



# Pulse Library

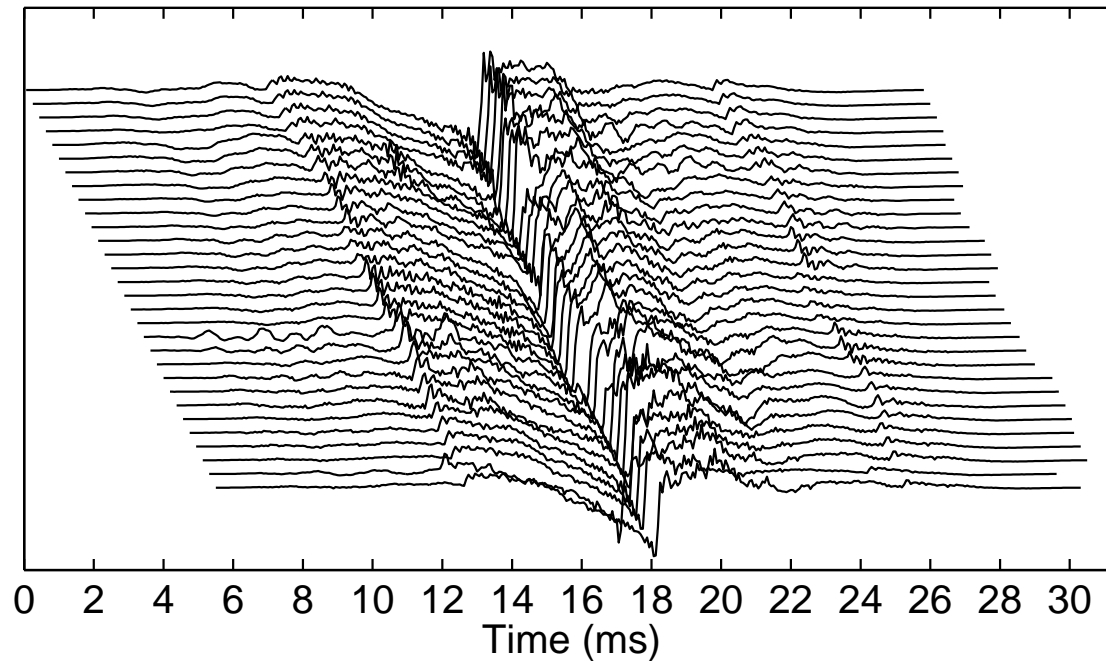
Construction of the pulse library:

1. Glottal closure instants (GCIs) are determined from the differentiated glottal flow signal
2. Each complete two-period glottal source segment is extracted and windowed with the Hann window
3. Pulses are linked with the corresponding voice source parameters:
  - Energy
  - Fundamental frequency (F0)
  - Voice source spectrum
  - Harmonic-to-noise ratio (HNR)
  - 10 first harmonic magnitudes

In addition, a down-sampled (10 ms) version of the pulse waveform is stored for evaluating concatenation cost in synthesis stage.

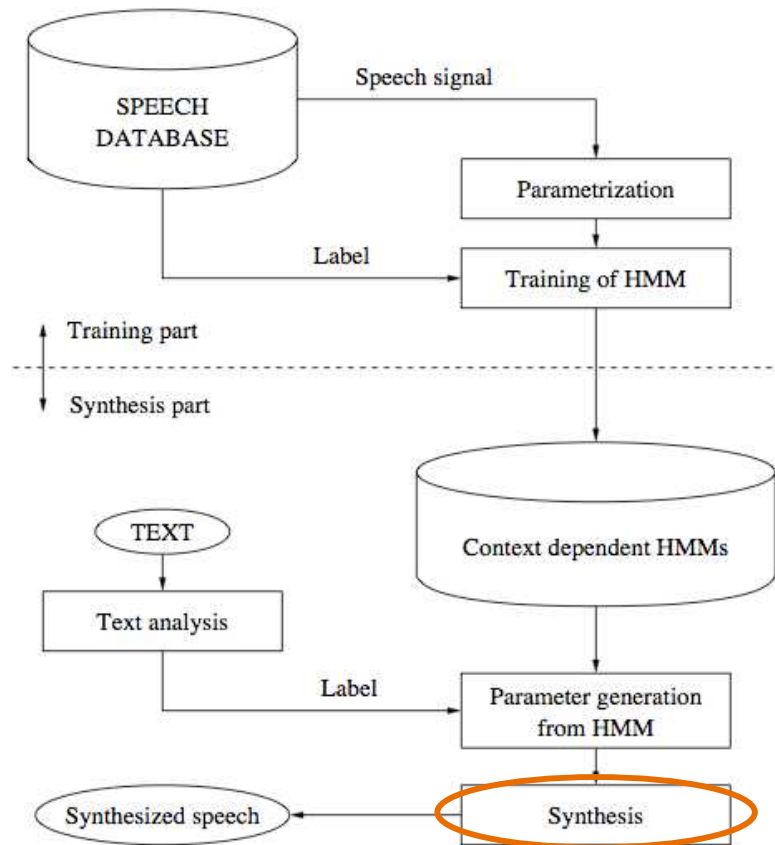
# Pulse Library

Consists of hundreds or thousands of glottal flow pulses (and the corresponding voice source parameters)



Windowed glottal volume velocity pulse derivatives from the pulse library of a male speaker

# Synthesis



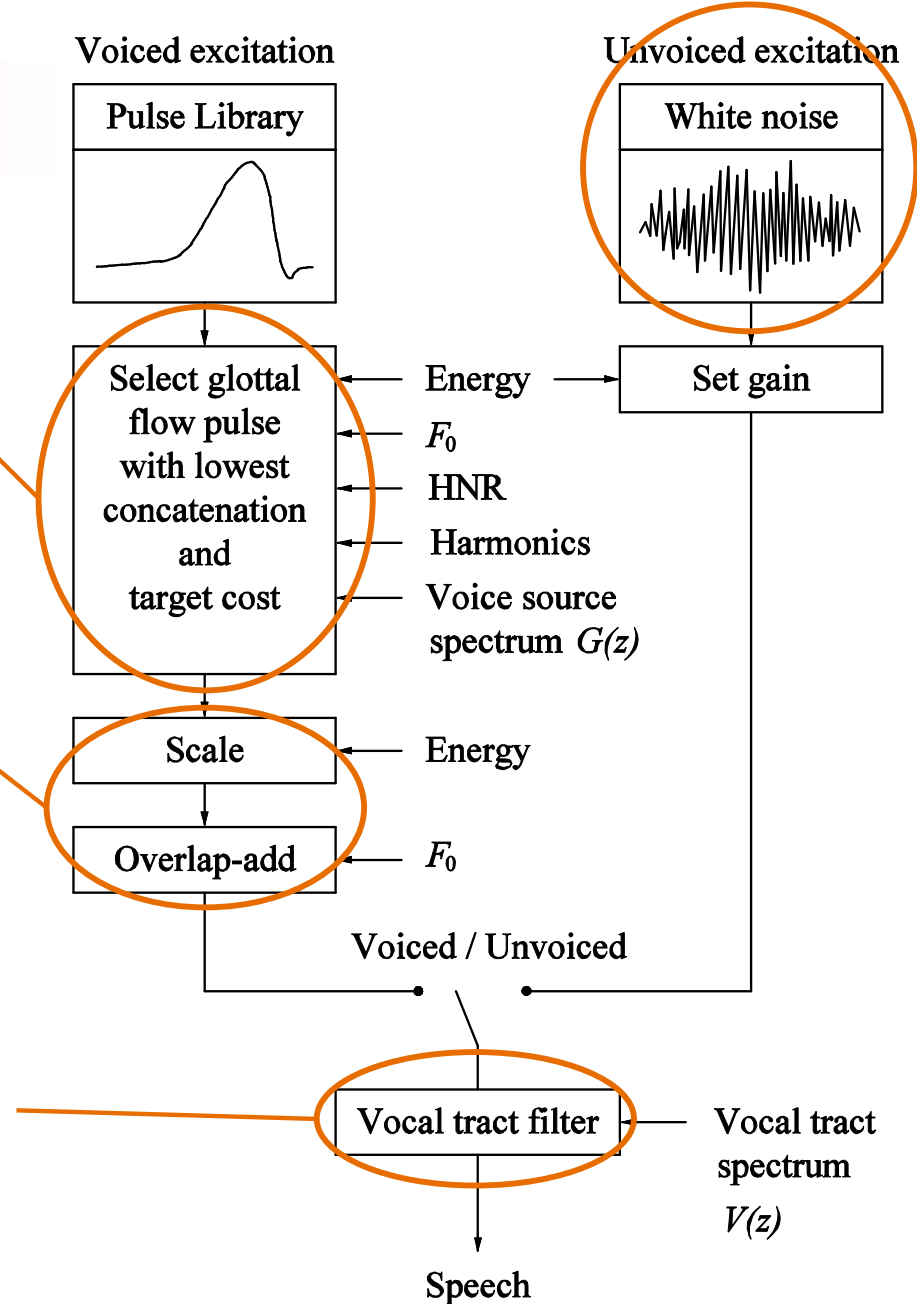
# Synthesis

In synthesis stage, excitation signal is generated by selecting the best matching pulses from the library according to the source features

Pulses are modified by scaling the magnitude and then overlap-added

White noise is used as unvoiced excitation

Finally, excitation is filtered with the vocal tract filter to generate speech



# Synthesis

The best pulse for each time index is selected by minimizing the joint cost composed of **target** and **concatenation costs**:

**Target cost:** RMS error between the voice source parameters generated by the HMM and the ones stored for each pulse.

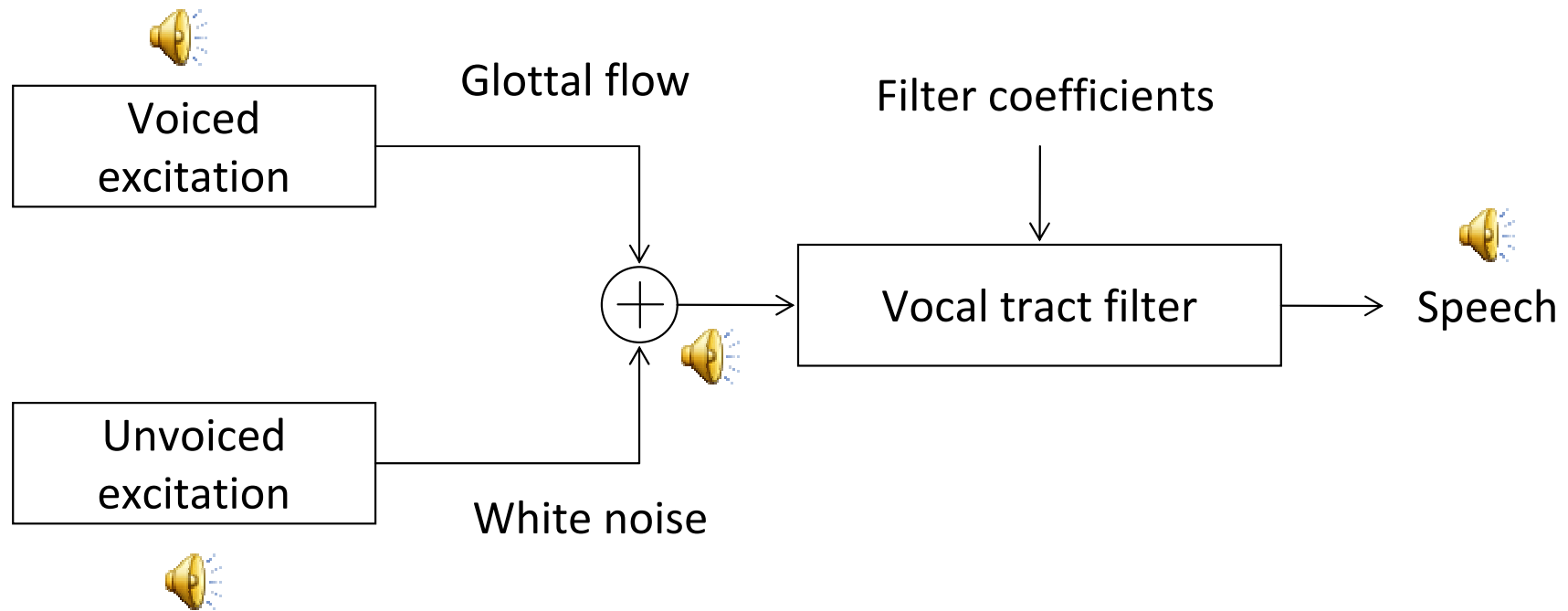
**Concatenation cost:** RMS error between the down-sampled versions of the pulse candidates

→ Scale energy of the pulse

→ Overlap-add

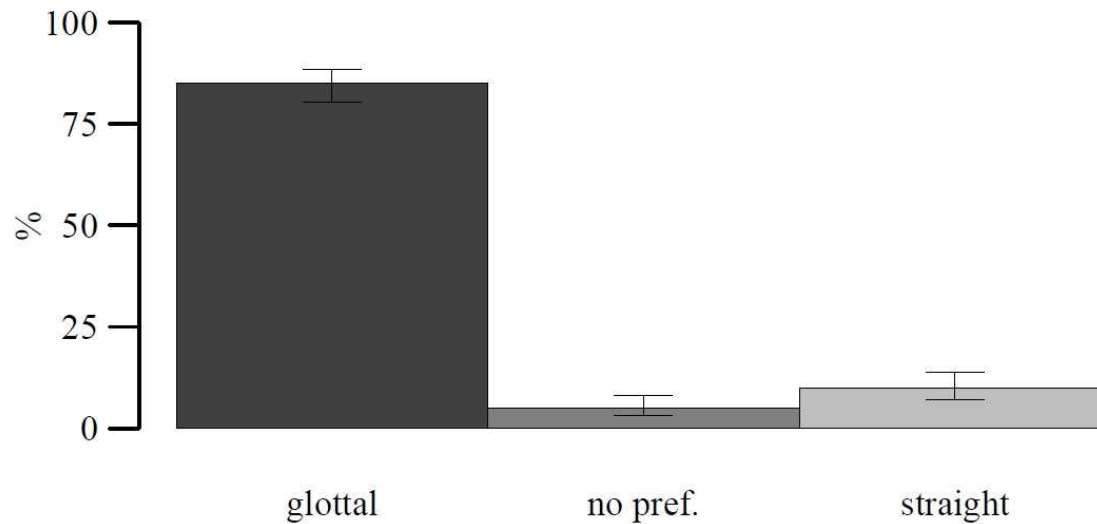


# Synthesis



## IV. Results and Samples

Previously, we have used only **one glottal pulse per utterance**.













Results of the listening test [2] comparing our synthesis method to the most widely used high-quality vocoder STRAIGHT.











# Single Pulse Technique

Samples:

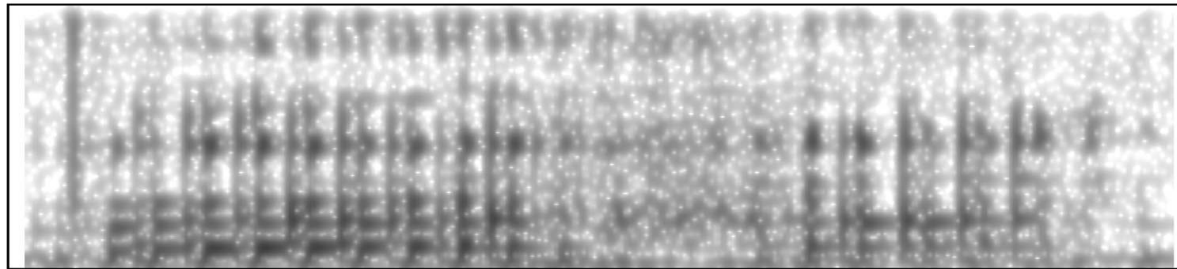
English	Male	Female
		
		

Blizzard Challenge 2010			
English			
Mandarin			

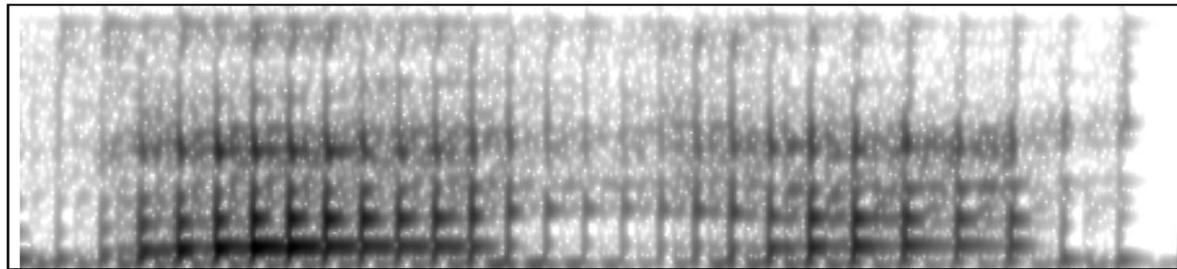
# Pulse Library Technique

Pulse library (ICASSP'11)	1pulse	pulselib
Finnish		
Finnish		
English		
English		

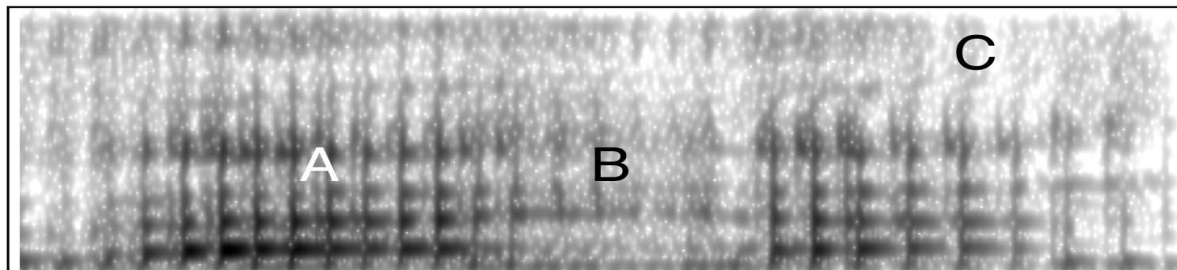
# Pulse Library Technique



Natural



Single pulse



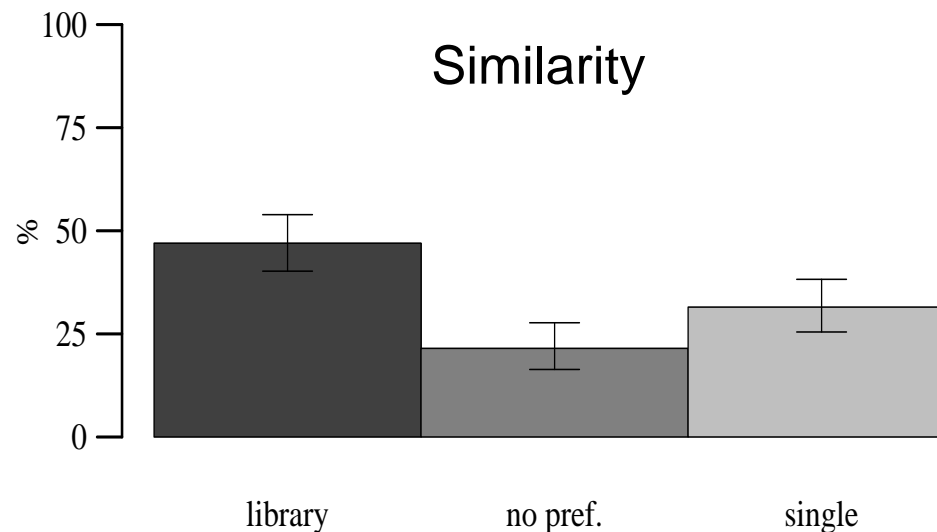
Pulse library

Spectrograms (0–8000 Hz) of the word “vähän” (little). Note the improved modeling of A) diphthony B) voiced fricatives C) high frequencies.

# Pulse Library vs. Single Pulse Technique

According to listening tests:

- ❑ Pulse library method is slightly preferred over the single pulse technique
- ❑ Better speakers similarity but creates some artifacts as well



- ❑ New physiologically motivated high-quality speech synthesizer
- ❑ Allows for better reproduction and control over the speech characteristics
- ❑ Pulse library generates more natural excitation and is preferred over single pulse technique

## References

- [1] P. Alku, “Glottal wave analysis with pitch synchronous iterative adaptive inverse filtering,” *Speech Communication*, vol. 11, no. 2–3, pp. 109–118, 1992.
- [2] T. Raitio, A. Suni, J. Yamagishi, H. Pulakka, J. Nurminen, M. Vainio, and P. Alku, “HMM-based speech synthesis utilizing glottal inverse filtering,” *IEEE Trans. on Audio, Speech, and Lang. Proc.*, vol. 19, no. 1, pp. 153–165, Jan. 2011.

Thank You!