



RAMCESS'07

**Realtime & Accurate Musical Control of
Expression in Singing Synthesis**

Table of Contents



Challenges

Previous Work

This Year's Aims

Database

Analysis

Synthesis

Difficulties in Synthesis

Challenges



Voice quality analysis/modification for expressive/emotional speech synthesis.

- 📌 Analysis: High quality source-tract decomposition is difficult : a long time challenge
- 📌 Synthesis: Control of low-level source parameters to achieve soft, tense, harsh... voices : research at its early stages

Real-time control for expressive singing synthesis.

- 📌 Design of appropriate UI and mapping functions

Previous work



GRAMCESS 1.0 :

Real-time

expressive singing

synthesis

Model-based voice

quality modification

/owel-only

synthesis

QuickTime™ and a
mpeg4 decompressor
are needed to see this picture.

This Year's Aims



include coarticulation

improve naturalness

support new voices : Able to synthesize a
specific person's voice




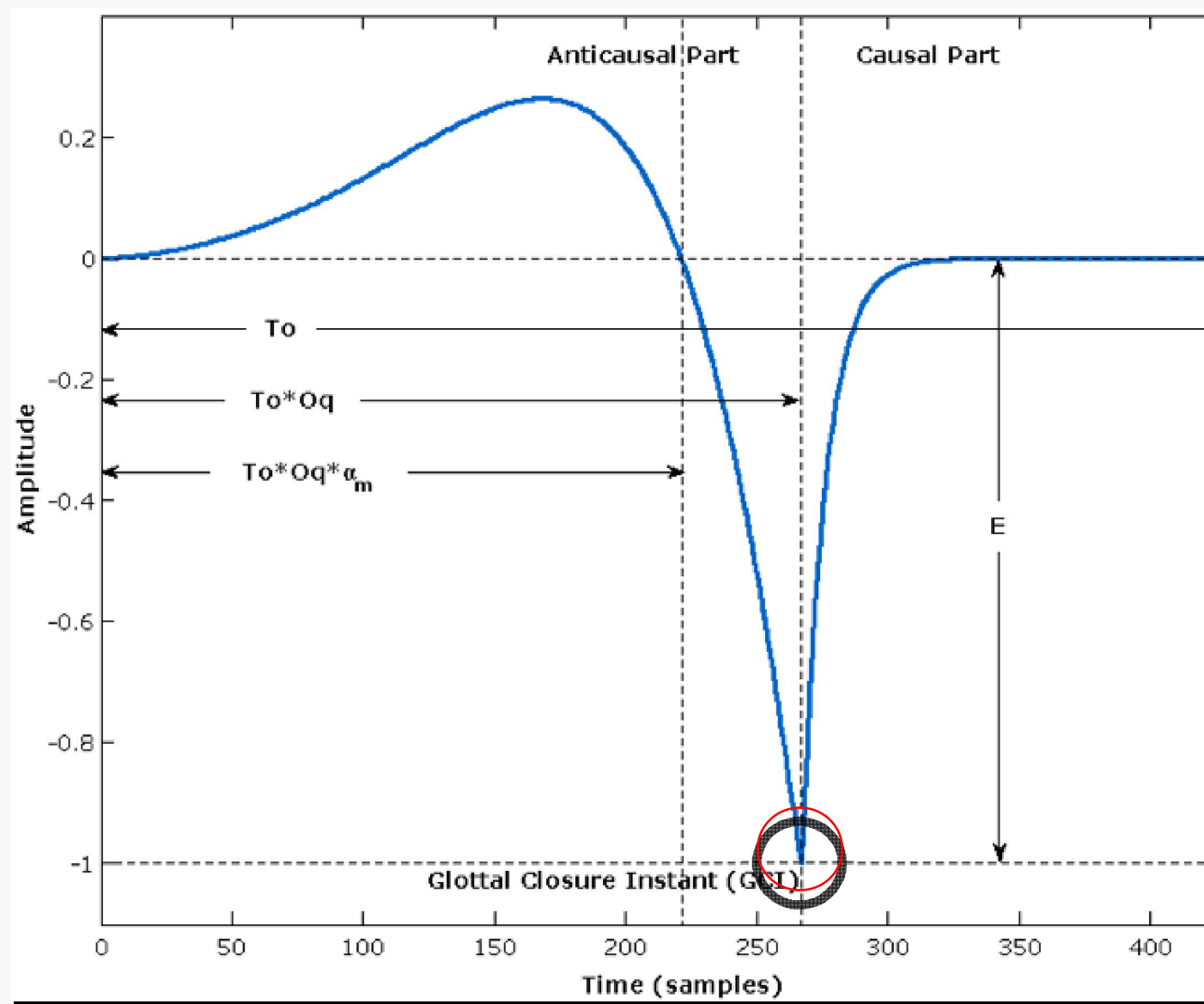
- Small size : ~ 2 seconds of recording
- Includes unvoiced consonants, vowels and transitions
- Constant pitch
- Constant voice quality (perceivably)



Background

LF Model : A model for glottal source

 **Glottal Closure Instant** : The main excitation point of and important feature for describing a



Background



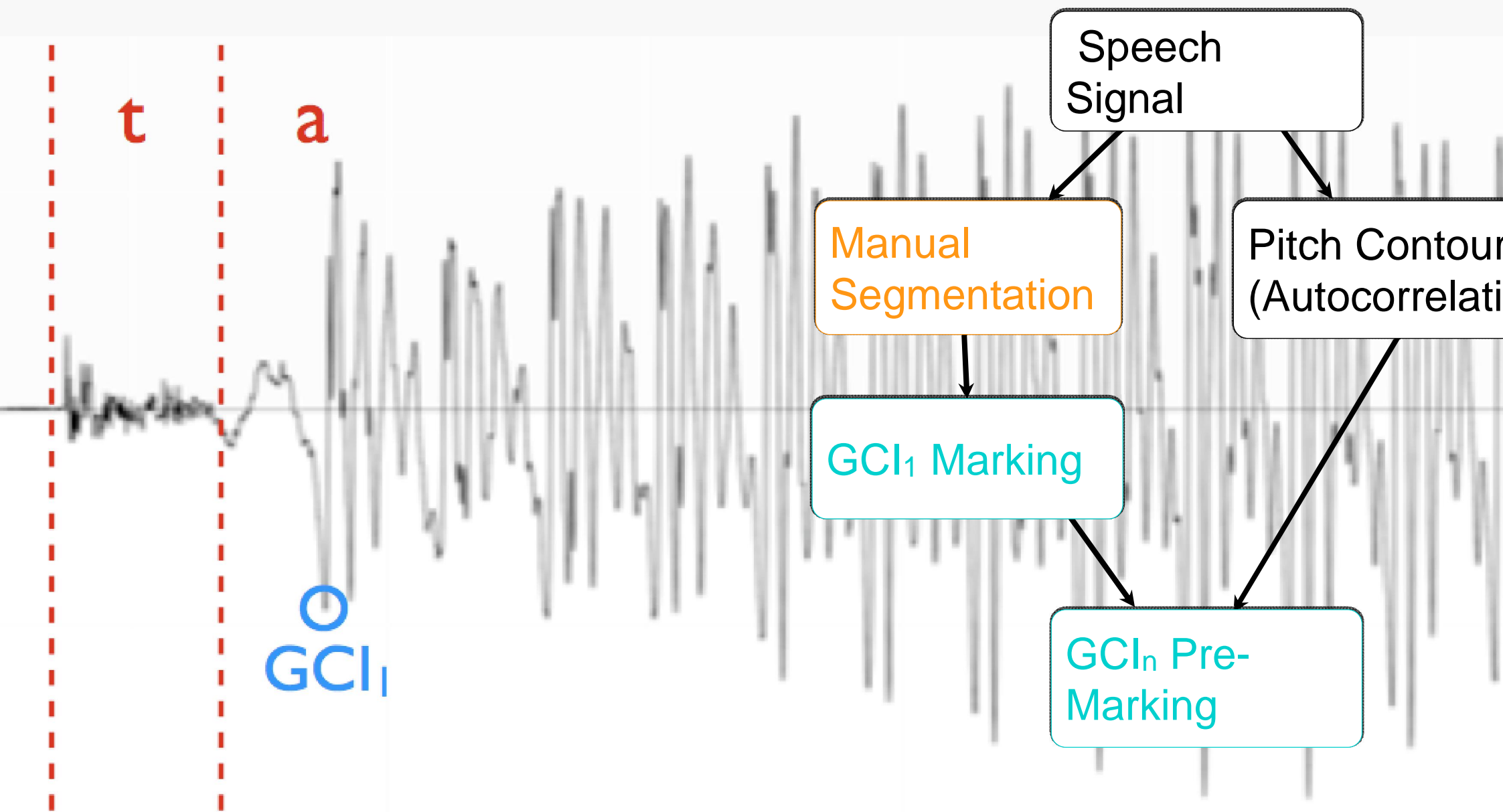
ZZT : Zeroes of Z-Transform

- 🔊 Source - tract decomposition
- 🔊 Sensitive to window position
- 🔊 High quality source estimation
- 🔊 Bad at spectral tilt

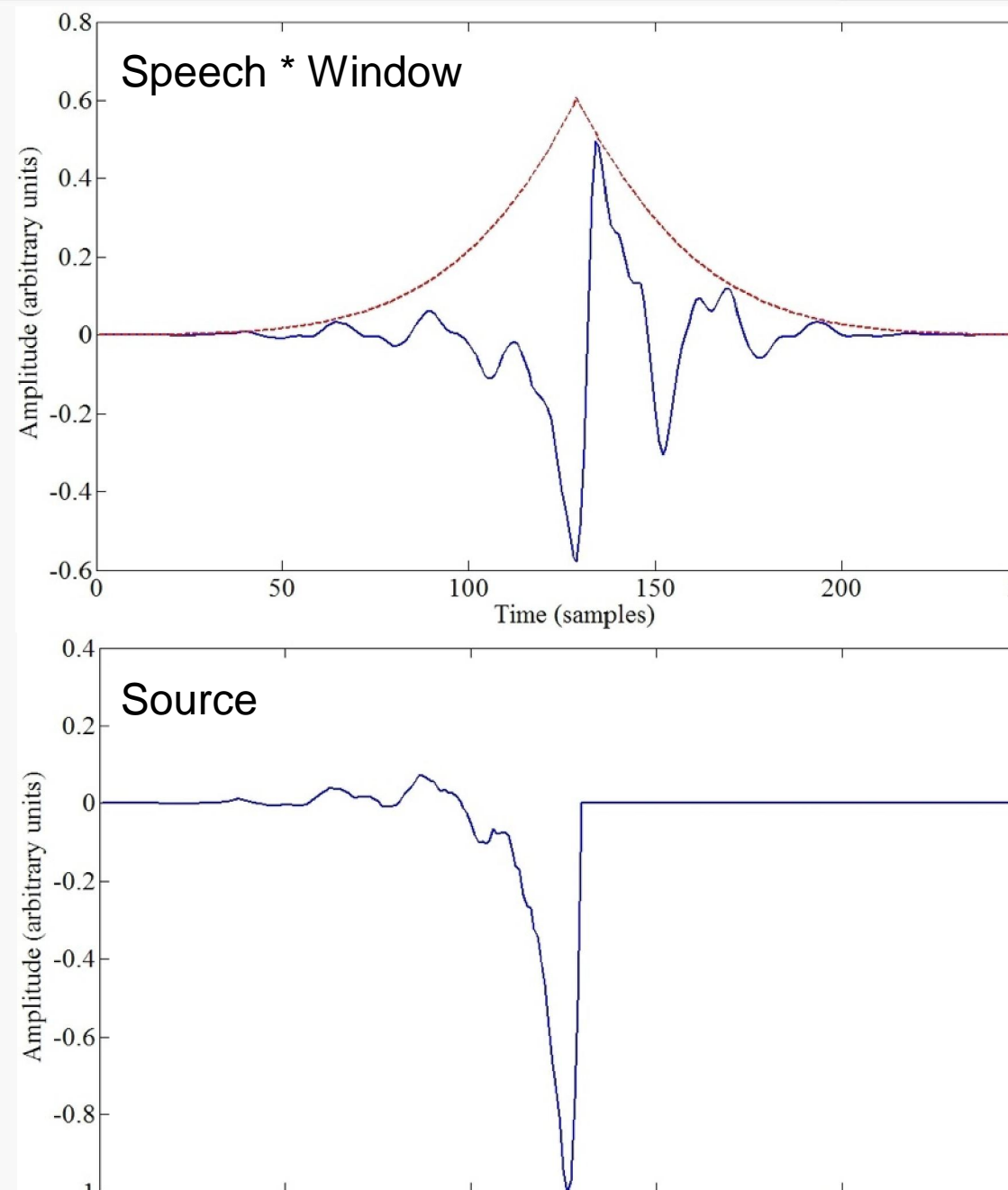
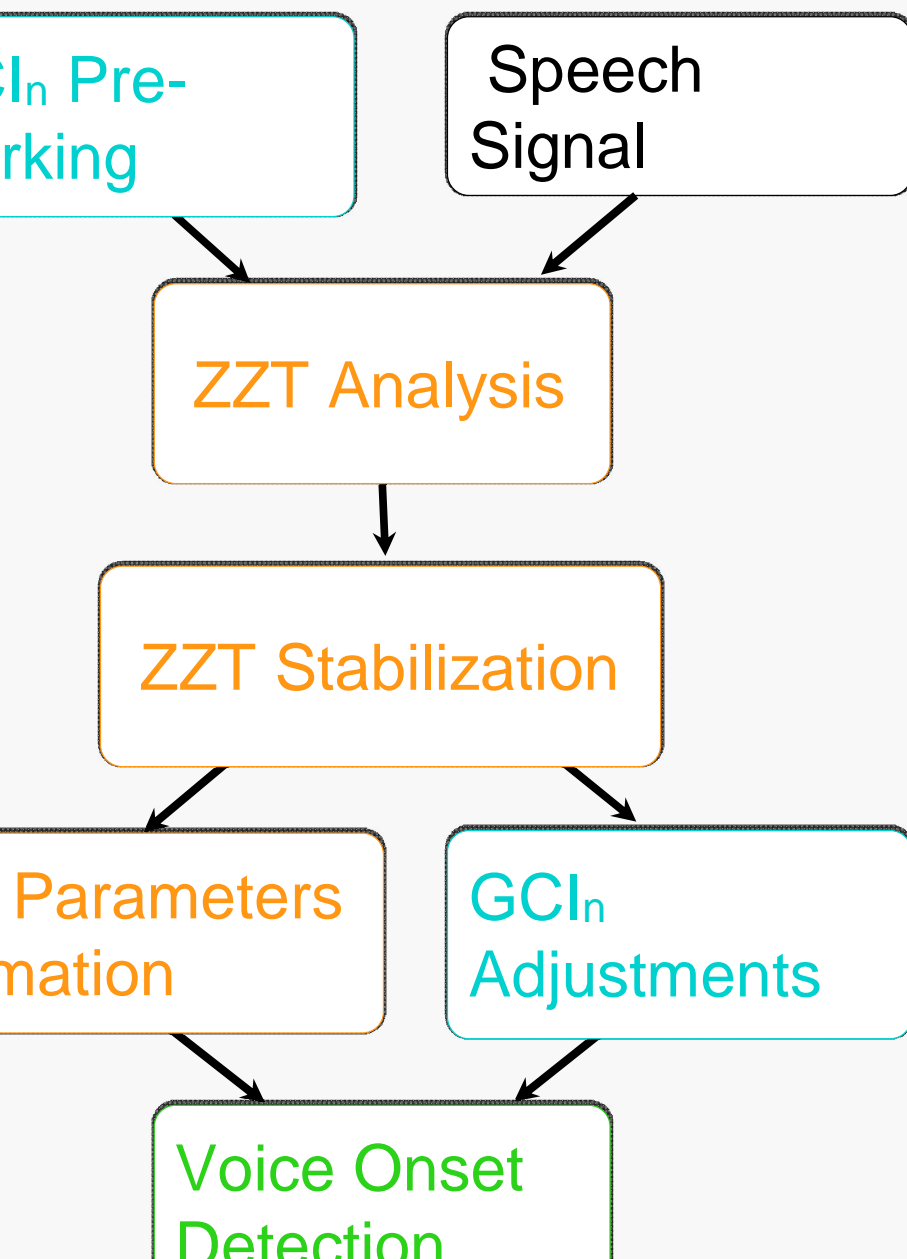
ARX : Auto-Regressive eXogenous

- 🔊 We use : An ARX model excited by the LF mode
- 🔊 Good at spectral tilt

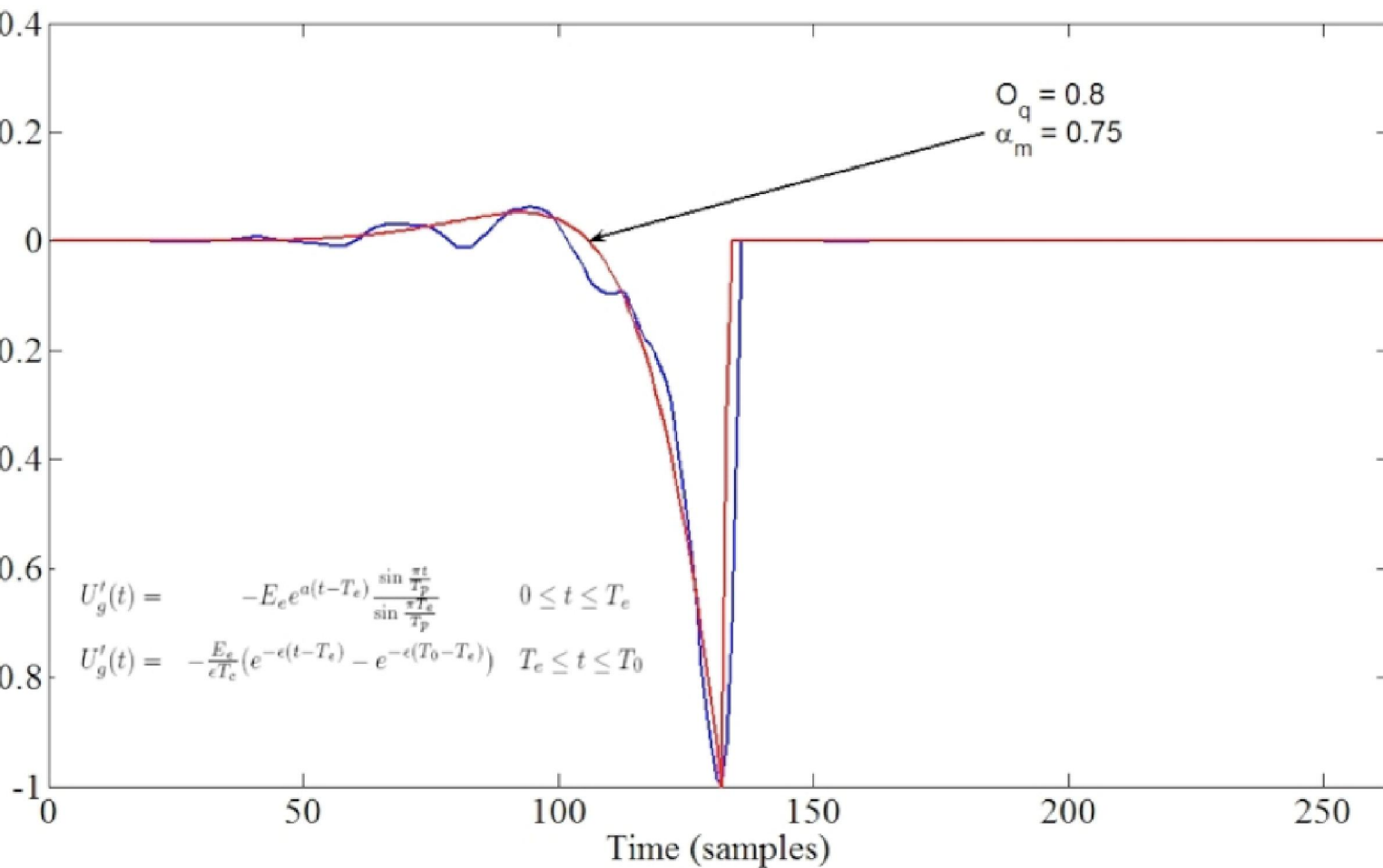
GCI & ALF Estimation



GCI & ALF Estimation



LF Model to AC part

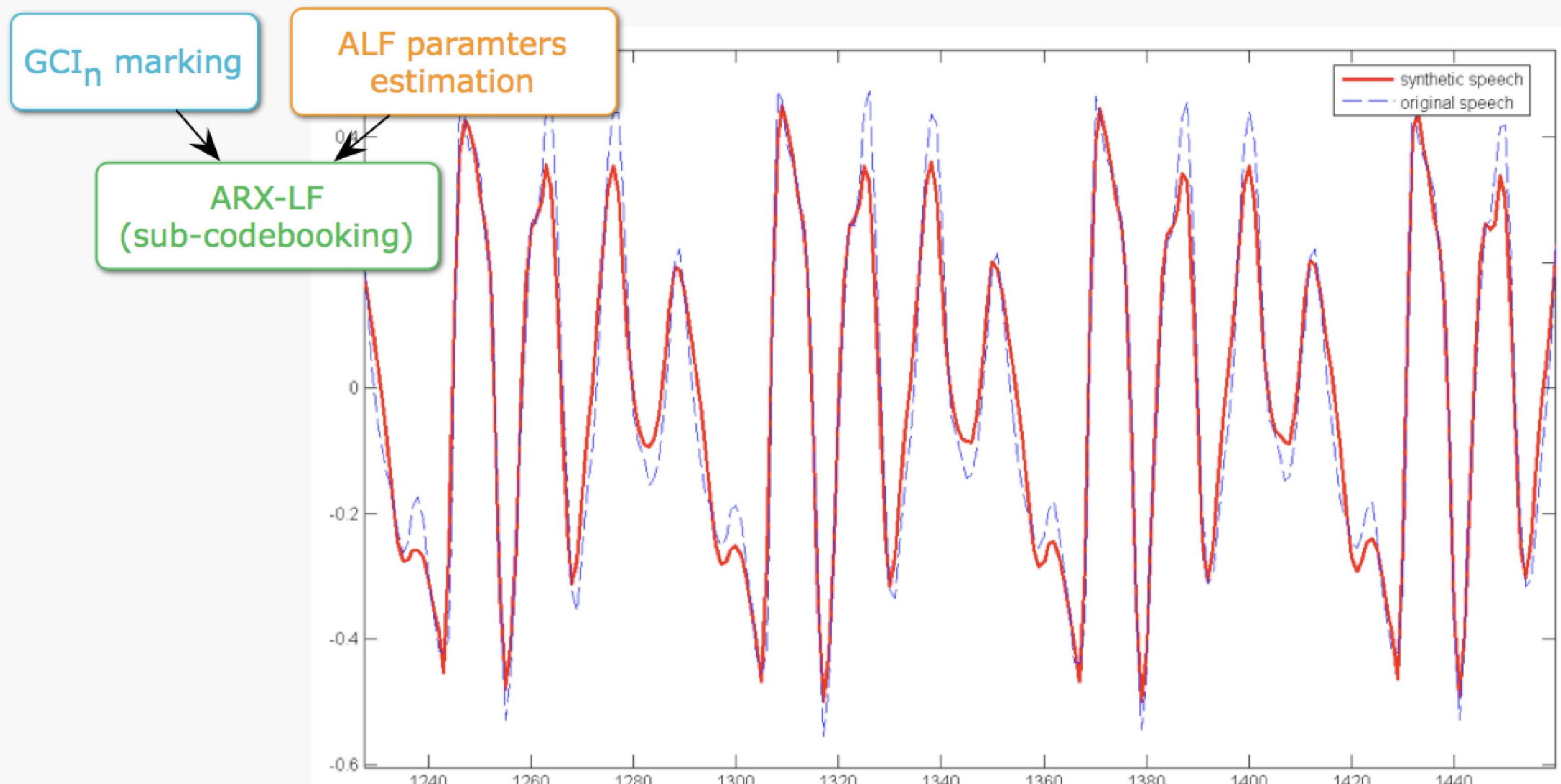


- Achieved by codebook search
- Performance on window signals
- Outputs

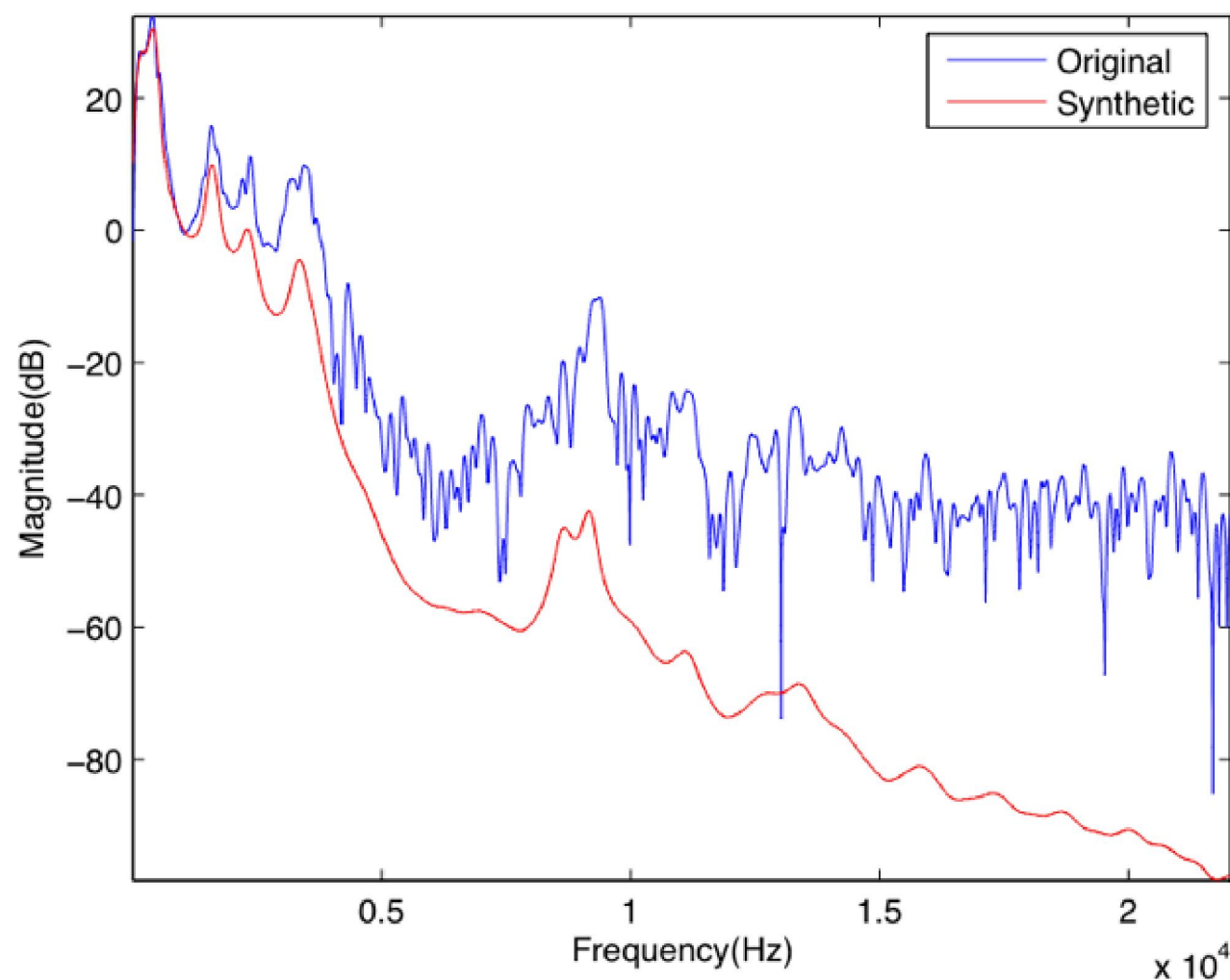
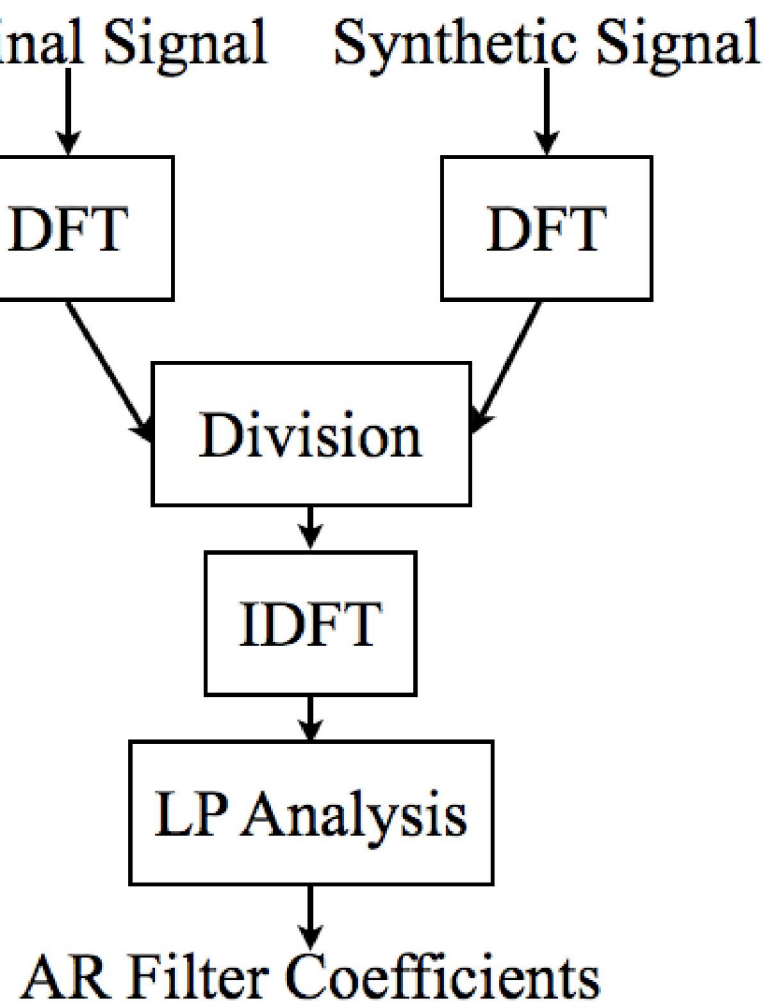
ARX for filter estimation



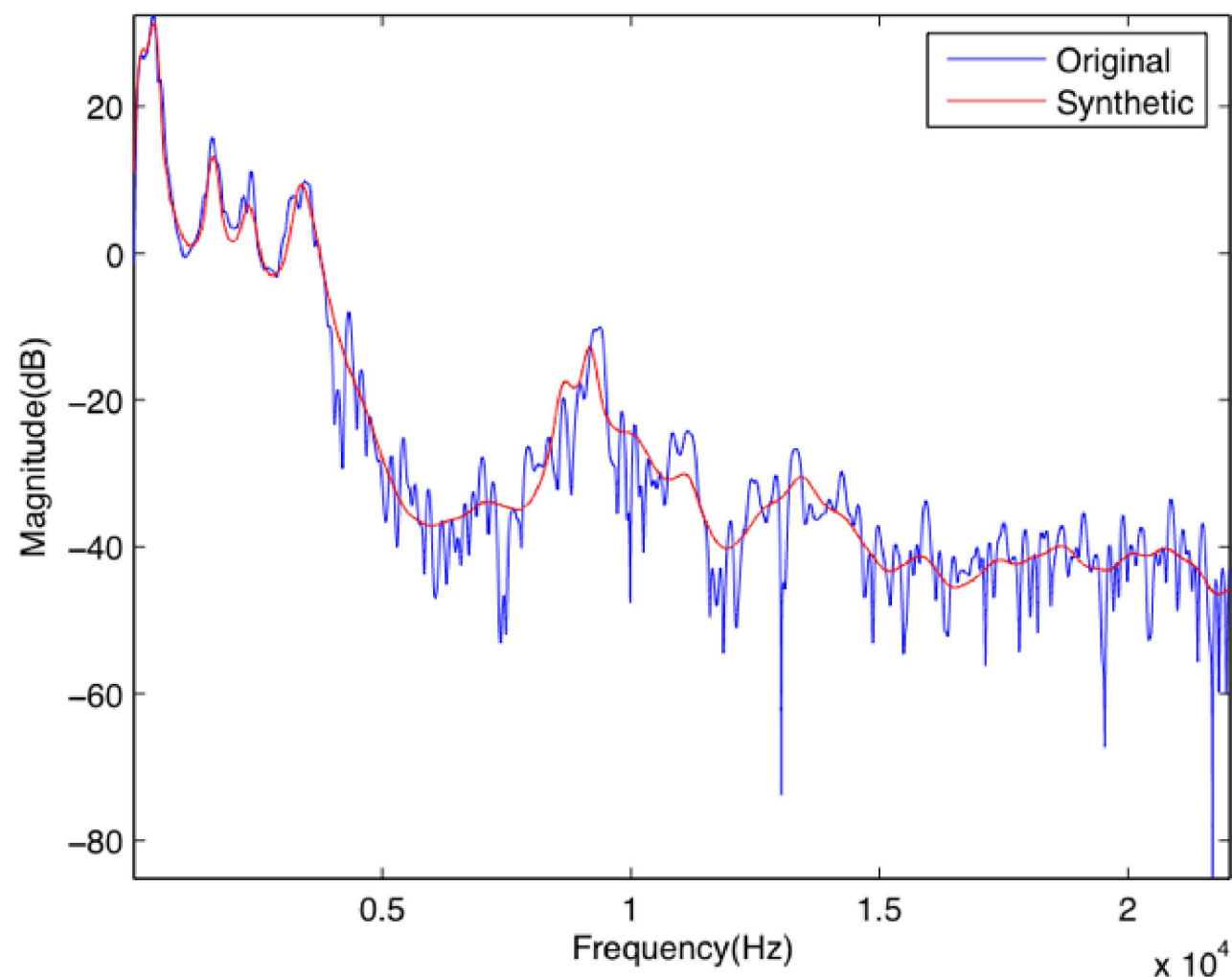
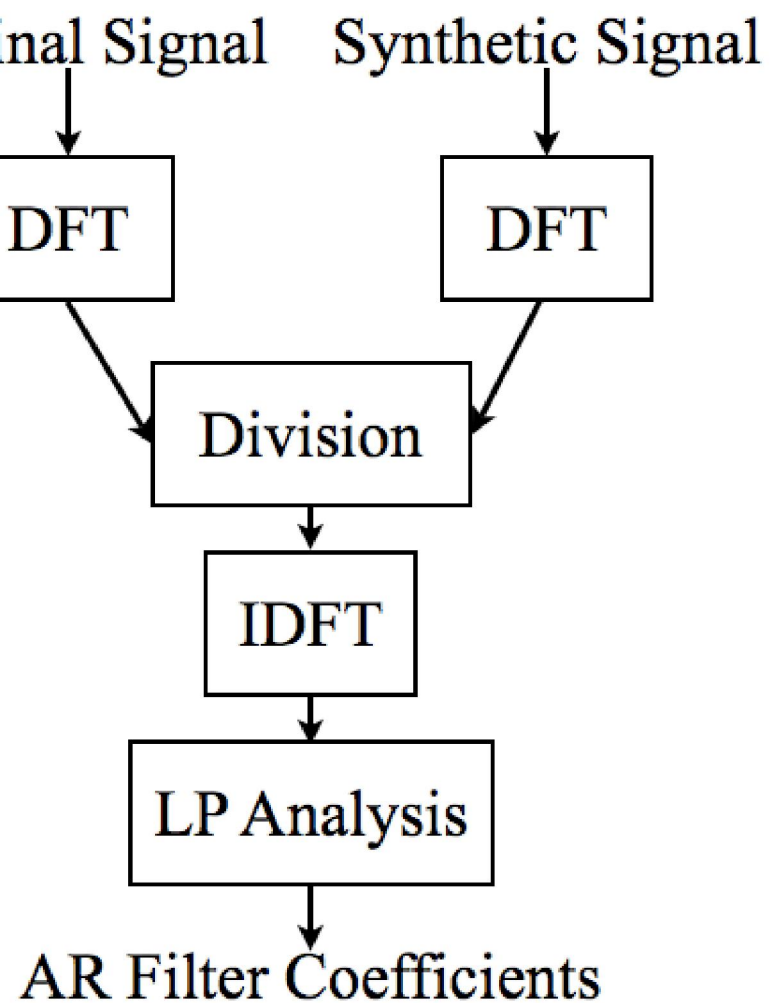
Joint estimation method for source and tract

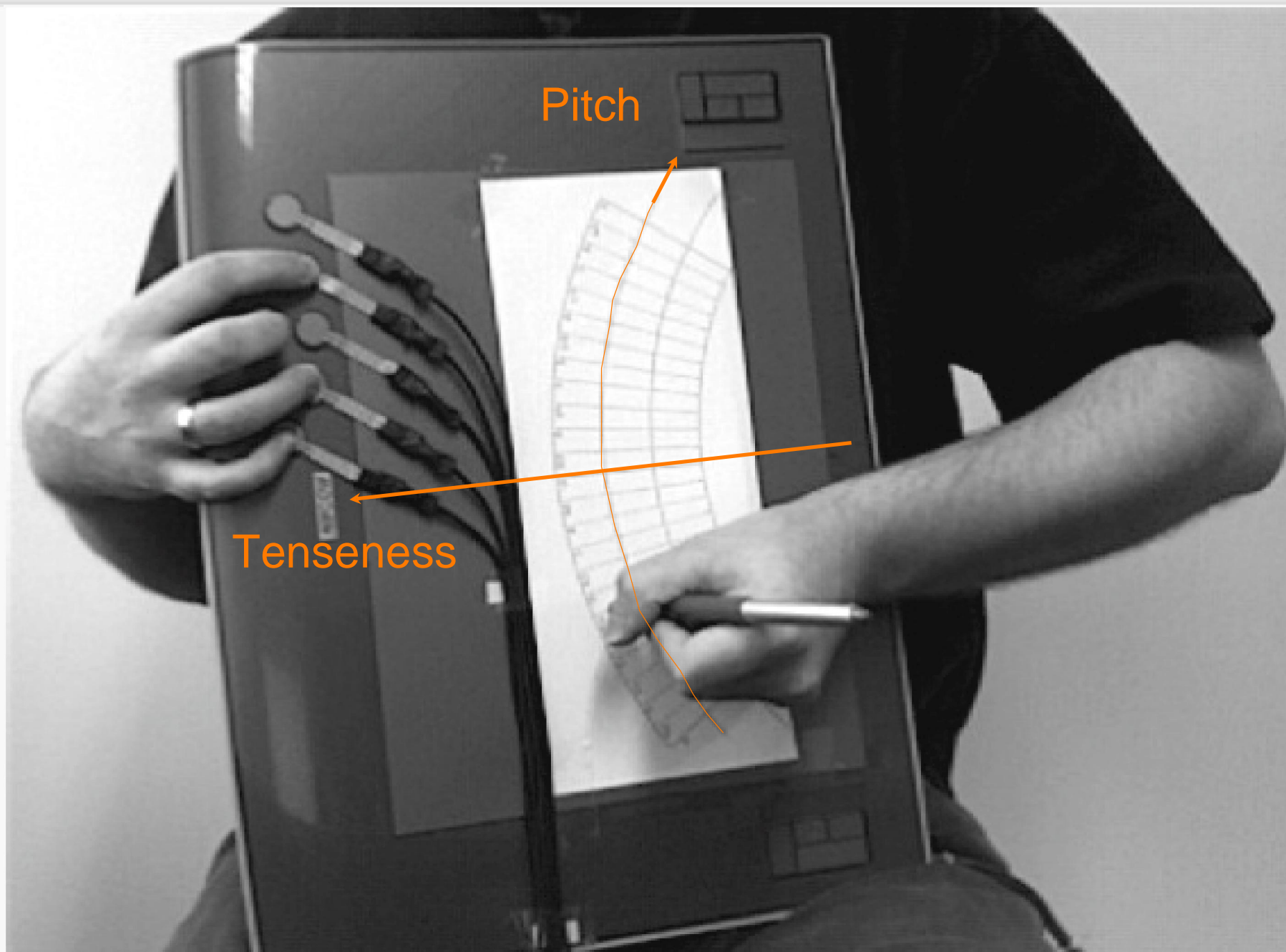


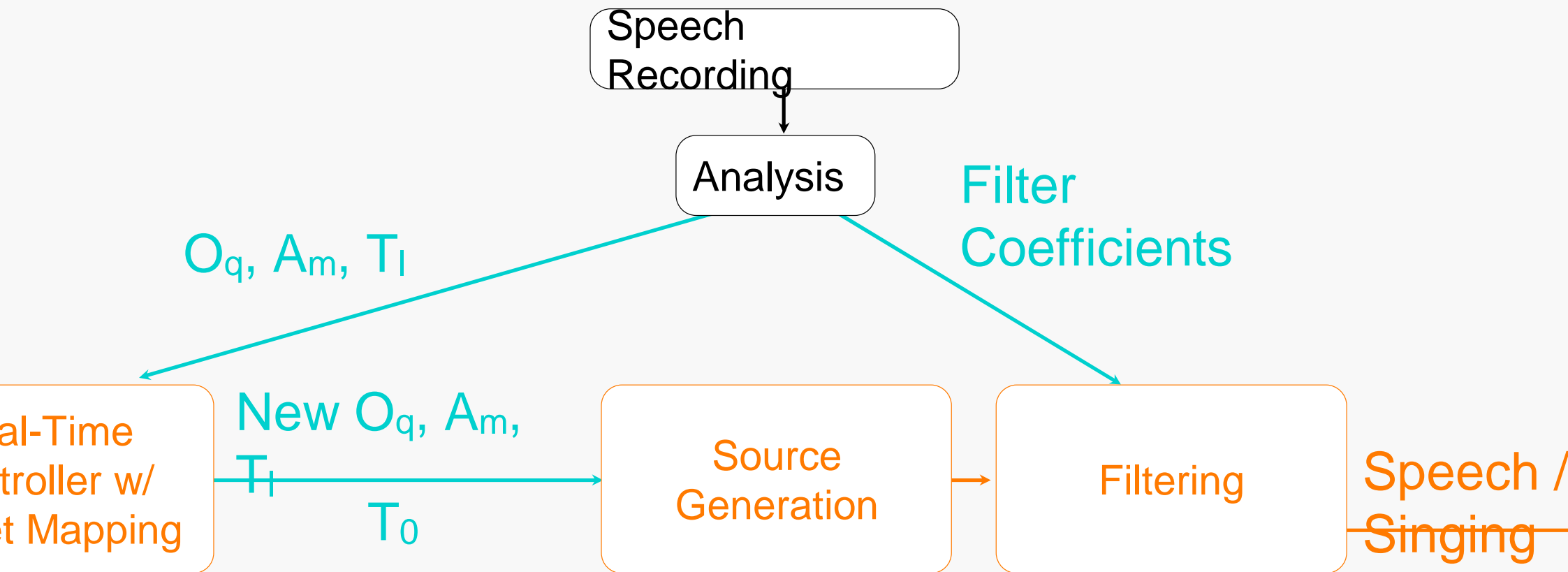
Compensation



Compensation









QuickTime™ and a
Photo - JPEG decompressor
are needed to see this picture.

Difficulties in Synthesis



Frame-by-frame analysis-synthesis results in rough speech

- 📌 Smoothing applied

Mapping of tenseness level to source parameters

- 📌 Obtained by trial and error

Real-time residual generation and mixing

- 📌 A method is devised and tested offline with positive results. Not integrated into the real-time



Improve transitions

Automating analysis

Devising a control interface for syllable
switching

Testing with different voices

Testing in musical performance



Thank you for your attention.

Questions?