# Speech Segmentation using Raw Sound

Merle Horne

Dept. of Linguistics

Center for Language and Literature
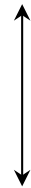
Lund University

Speech synthesis/Speech recognition/Speech understanding

$\updownarrow$

Psycholinguistic models of speech production/comprehension

# The role of function words in spontaneous
# speech processing

**MERLE HORNE**
**JOHAN FRID**
**GÖSTA BRUCE**
**BIRGITTA LASTOW**
**MIKAEL ROLL**
**ADINA SVENSSON**

www.ling.lu.se/projects/ProSeg2.html

.

# Psycholinguistic hypotheses tested on spontaneous speech

- **Commit- and- Restore hypothesis**: Stranded function word reflect "syntactic commitments" (Clark & Wasow (1998)). I.e. they signal that the speaker intends to produce a constituent of the kind signalled by the function word produced

- **Complexity hypothesis**: the probability that a speaker will hesitate in speech production will increase, the more complex the constituent being planned is (Clark & Wasow (1998))

Related Neurophysiological hypotheses (Pulvermüller 1995, 2003)

- Function words stored in the perisylvian cortex (Broca's region) - lateralized

  - Content words have a more widely distributed storage - not lateralized- support from imaging studies (Pulvermüller 1995, 2003)
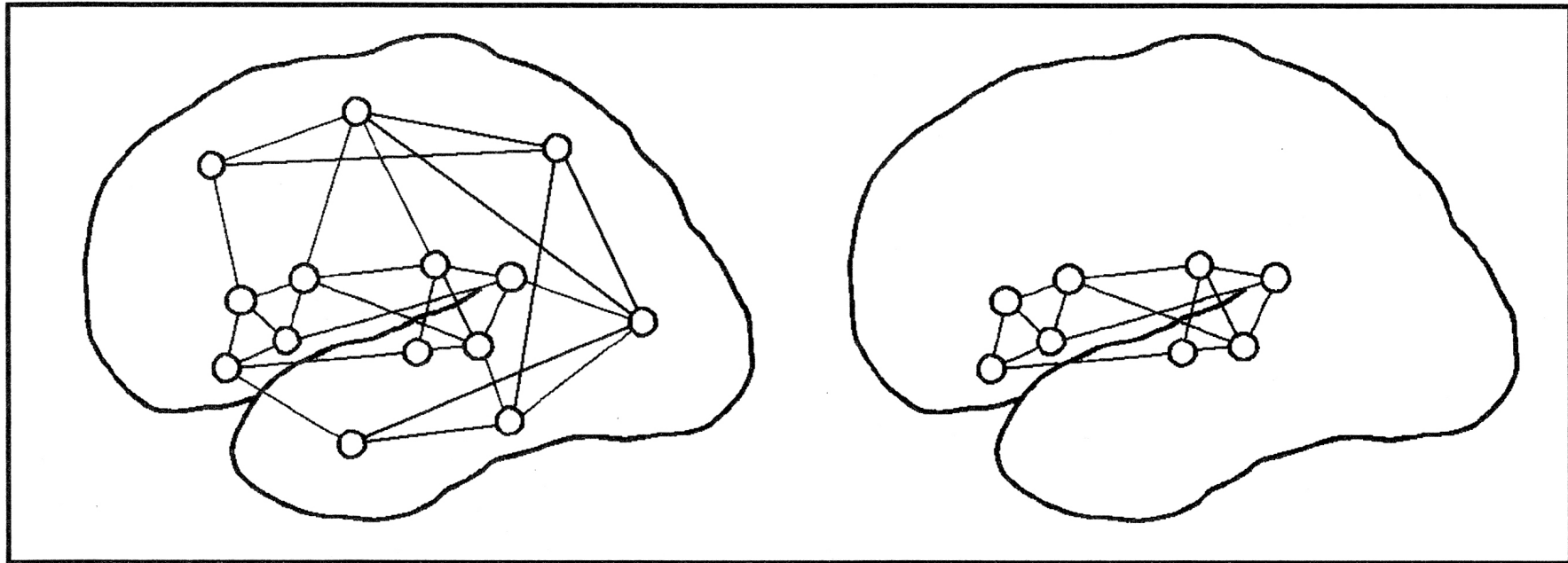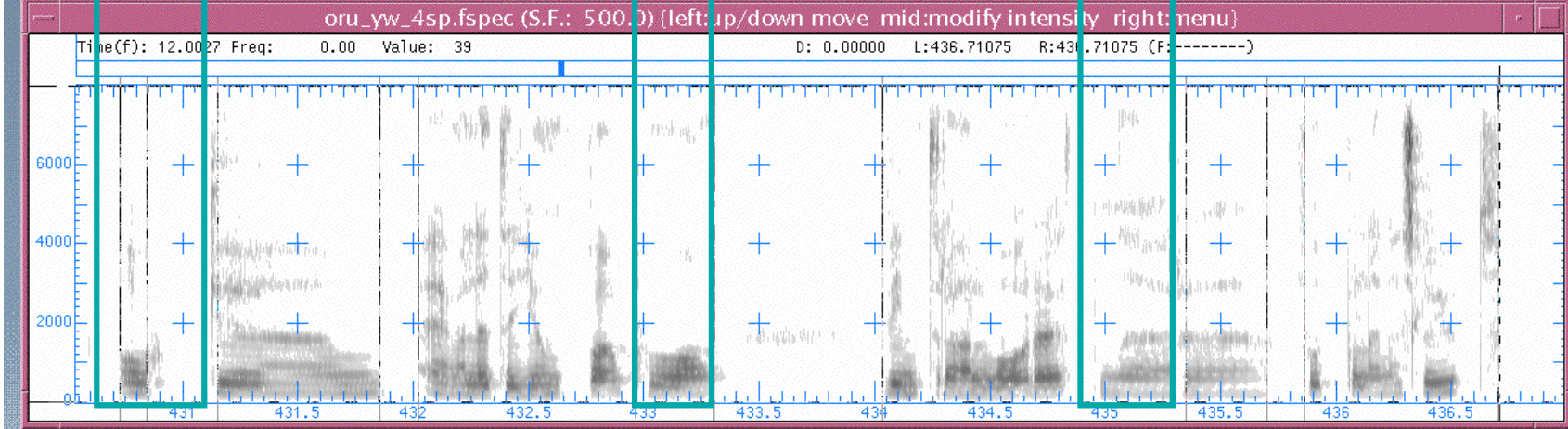
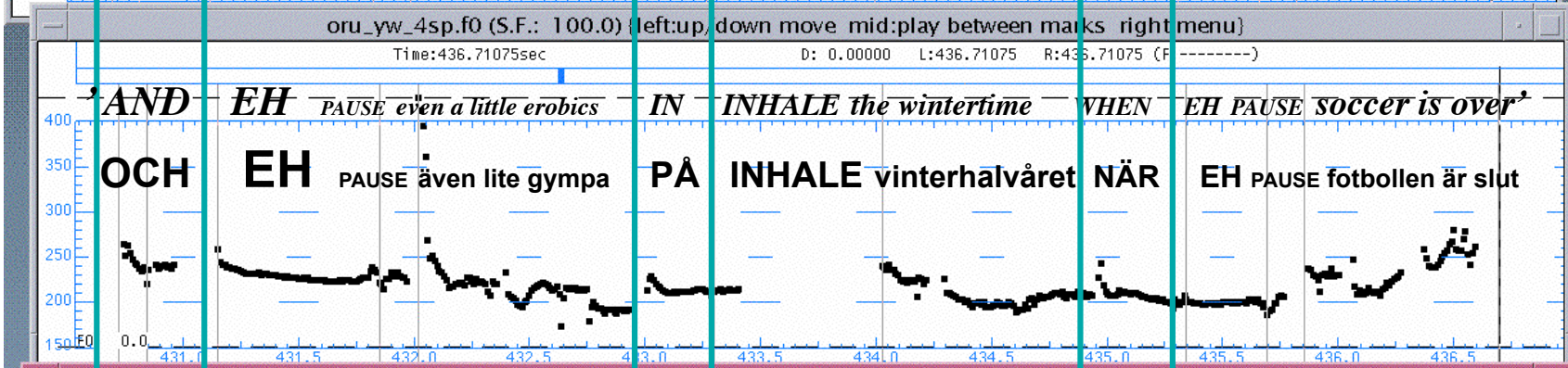**content word**             **function word**

**Figure 6.7.** Schematic illustration of left-hemispheric distributions postulated for cell assemblies representing high-imageability content words (left) and highly abstract function words (right).

(From Pulvermüller 2003)

# Speech 'Chunking' problem (Miller 1956)

Units for speech processing

(speech synthesis/speech recognition)

Timing restrictions on speech processing

- Memory research (Baddley 1997)
  'Phonological loop' (2 sec.)
- Cognitive linguistics (Chafe 1994)
  'Focus of consciousness' (1-2 sec)

CONCEPTUALIZER
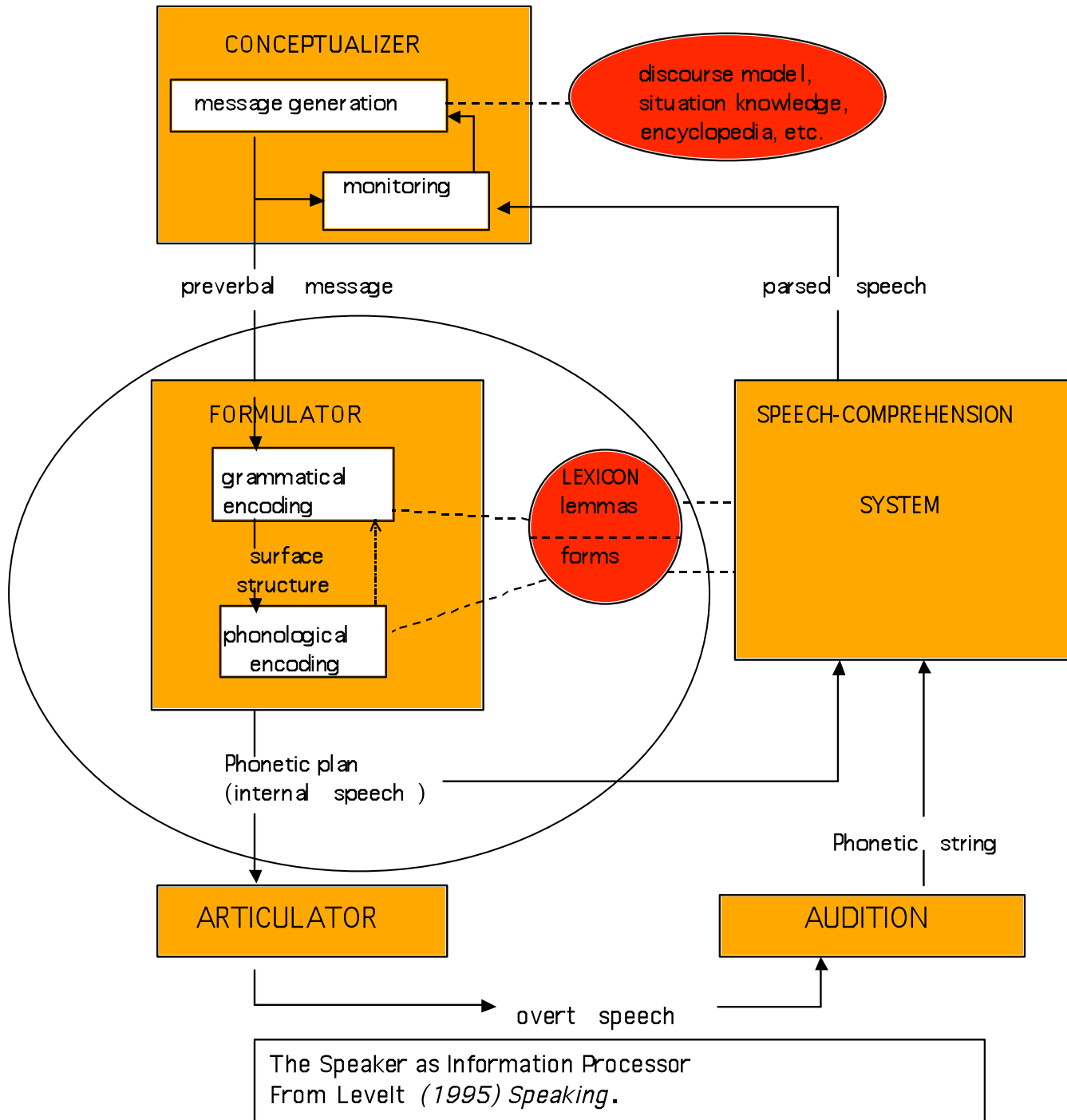
message generation

monitoring

discourse model,
situation knowledge,
encyclopedia, etc.

preverbal message

parsed speech

Speech planning unit

Ca. 2-2.5 sec. time limit

FORMULATOR

grammatical encoding

surface structure

phonological encoding

LEXICON
lemmas

forms

SPEECH-COMPREHENSION

SYSTEM

Phonetic plan (internal speech)

ARTICULATOR

AUDITION

Phonetic string

overt speech

The Speaker as Information Processor
From Levelt (1995) Speaking.

too_ow_1sp.wav (S.F.:16000.0) {left:up/down move  mid:play between marks  right:menu}

Time:476.81463sec          D: 3.56169    L:474.94025    R:478.50194 (F:    0.28)

10000

0

-10000

samples    148

too_ow_1sp.f0 (S.F.: 100.0) {left:up/down move  mid:play between marks  right:menu}

Time:476.81463sec          D: 3.56169    L:474.94025    R:478.50194 (F:    0.28)

*'so that then we moved down the shop        there'*

så  att då flyttade vi ner  butiken        dit

400
350
300
250
200
150
100

400
325
250
175
100

Hz

475.5        476.0        476.5        477.0        477.5        478.0

too_ow_1sp

s | v
  | att

INHALE PAUS                    så att då flyttade vi ner butiken dit          P

INHALE PAUS   så | att |      flyttade | vi | ner       butiken       dit |   PAUS   INHALE
              <p2>  | då                                            </p2>

<t2>

Ca. 2.5 sec.

10000

0

samples

too_ow_1sp.f0 (S.F.: 100.0) {left:up/down move  mid:play between marks  right:menu}

'then they had to help very much with it' 'Since then we had to get up PAUSE redo INHALE'

400

325

då fick de ju hjälpa till väldigt mycket med det   för att då skulle vi ju ha upp PAUS göra om INHALE

250

H%

175

100

0.0

Hz

384.5    385.0    385.5    386.0    386.5    387.0    387.5    388.0    388.5    389.0

too_ow_1sp

s  o
att

ACK  PAUS                                                för att då skulle vi ju ha upp PAUS göra om    INHALE
INHALE   då fick de hjälpa till väldigt mycket med det

ACK PAUS  fick           ju        till        mycket      det   att  skulle   ju      upp   PAUS      om
INHALE  då    de    hjälpa    väldigt   med   <p2>  då      vi ha           göra    INHALE

2><t2>                                   <t2/><t2>                              </t2><t2>

Ca. 2 sec.                                      Ca. 2 sec.

One        Two

Time (s)

Pitch (Hz)

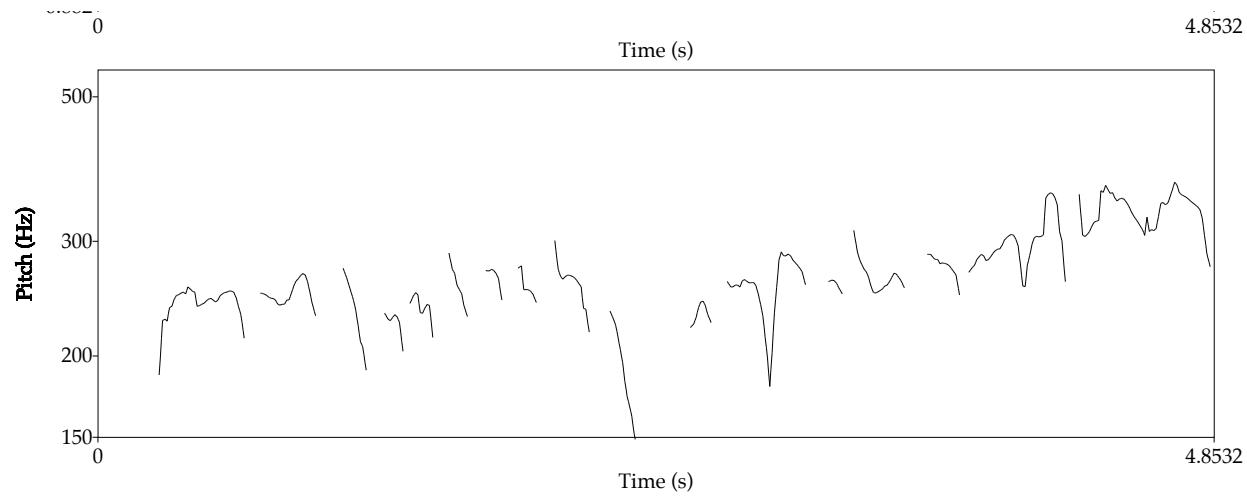Time (s)
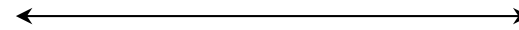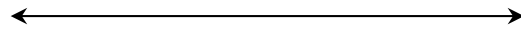
[Jenny just nu nästsist] [hon kan klippa Tchechenko]

[Å så är det hårt där framme] [å så kommer Johanna Hayes och ramlar]

Ca.2.1 sec

Ca. 2.3 sec

# Inhalations as anchors for speech segmentation

- Can one automatically recognize inhalations?
  - Characteristics: noise, lack of F0, can have formant structure if oral
- Possible method: Template matching
  - Uses distance measures between a reference sound and signal being processed

# Using cepstral coefficients for inhalation pause detection in spontaneous speech

*Anders Johansson, Johan Frid, Merle Horne*
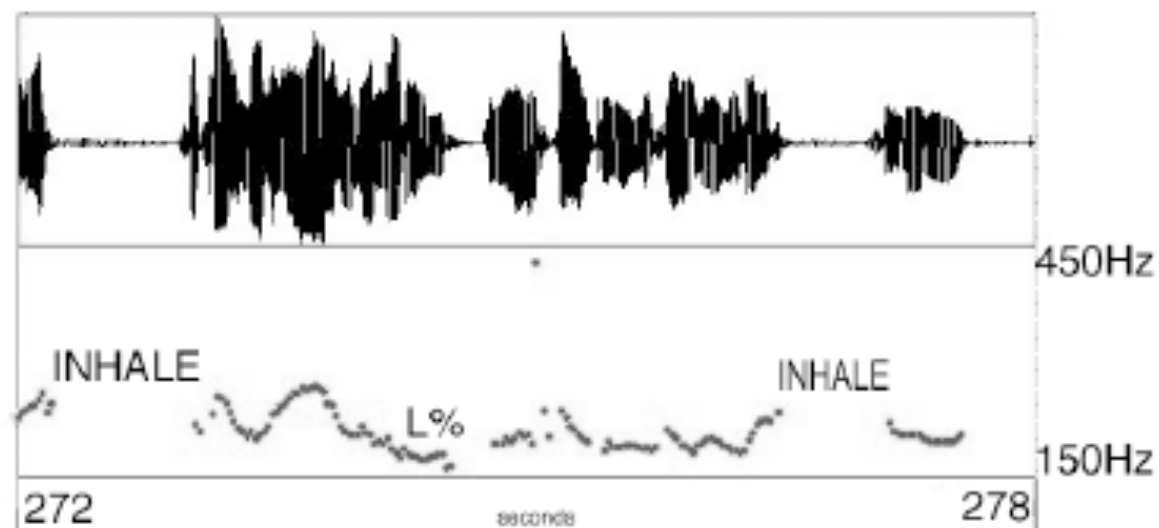
SPECOM 2005, Patras,
Greece



Figure 1: *An example of spontaneous speech illustrating how inhalations (labelled INHALE) can be used as anchors in the segmentation of speech into processing units. The speech between the inhalations consists of two clauses:* Så tränar man med jämna mellanrum *'So you train regularly' and* det gick pågick ju under flera månader *'it last lasted for several months'. The clauses constitute two prosodic phrases separated by a pause.*
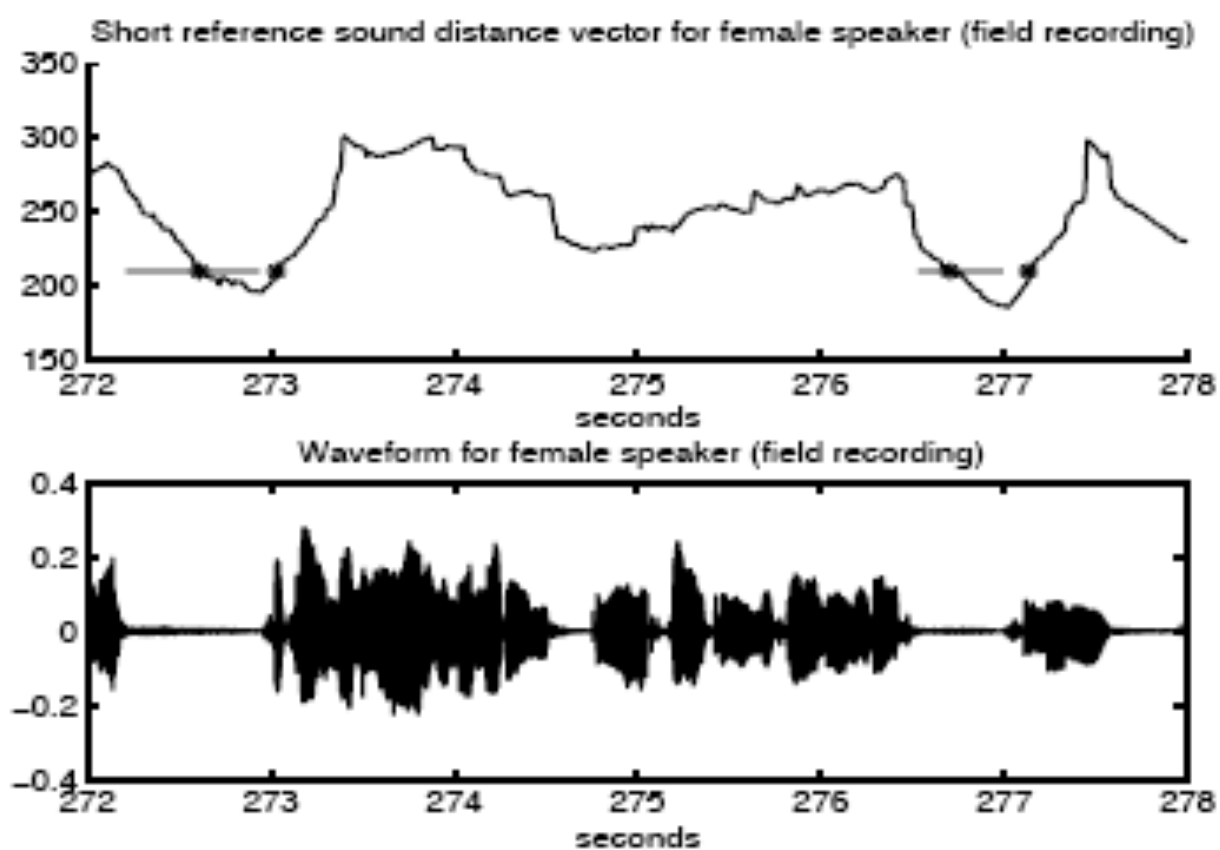
Figure 5: *Distance vector of field recording of female speaker as generated from short reference sound. Stars indicate detected inhalations, lines indicate tagged inhalations. Recording from the SweDia material.*
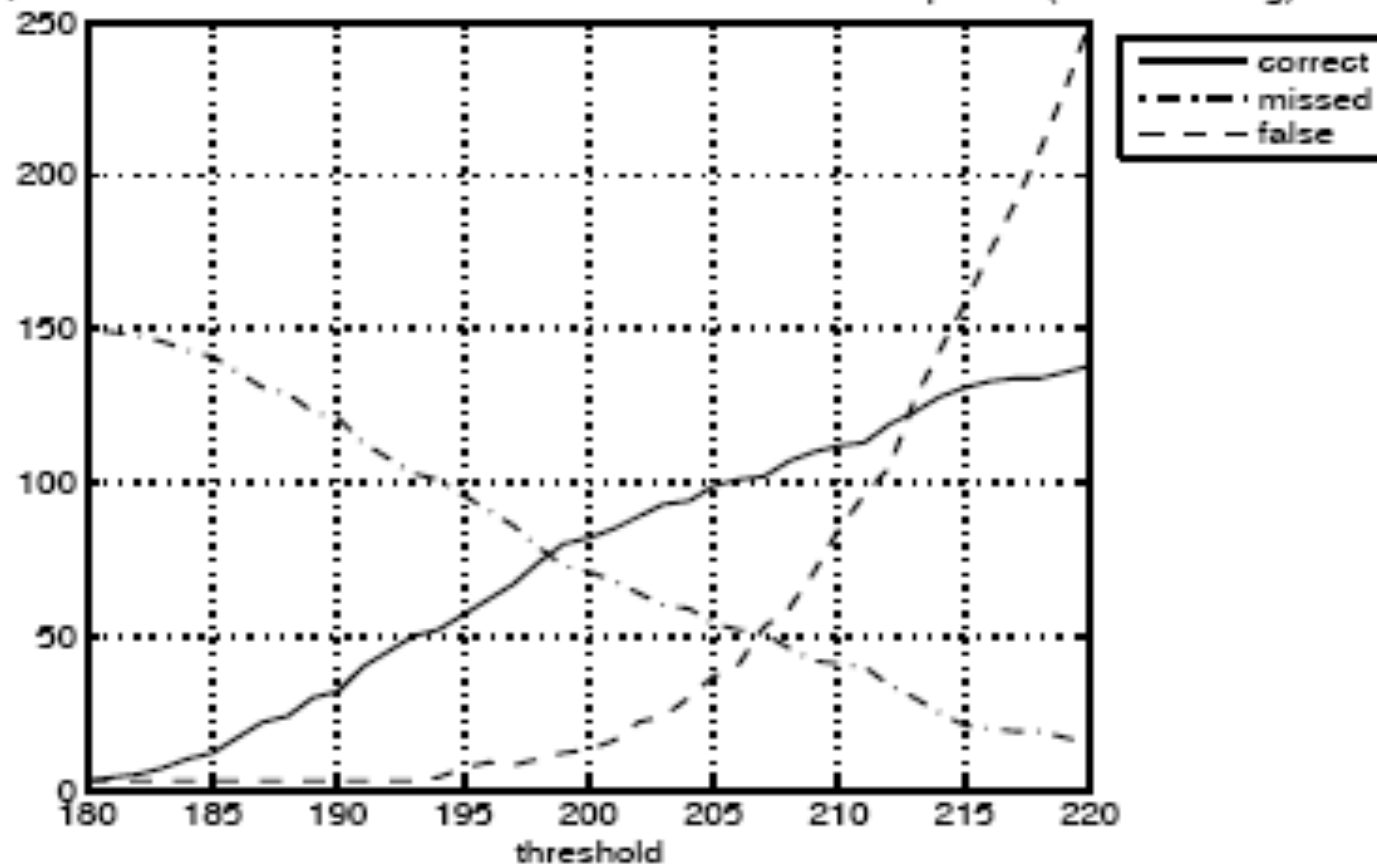
Figure 9: *Number of correct, false and missed identifications of inhalations in field recording of female speaker as generated from short reference sound.*

Observations:

False identifications associated with sounds of the following types: exhalations, word-final aspirated sounds, whispers and voiceless fricatives

Properties of inhalations quite stable, even over different speakers