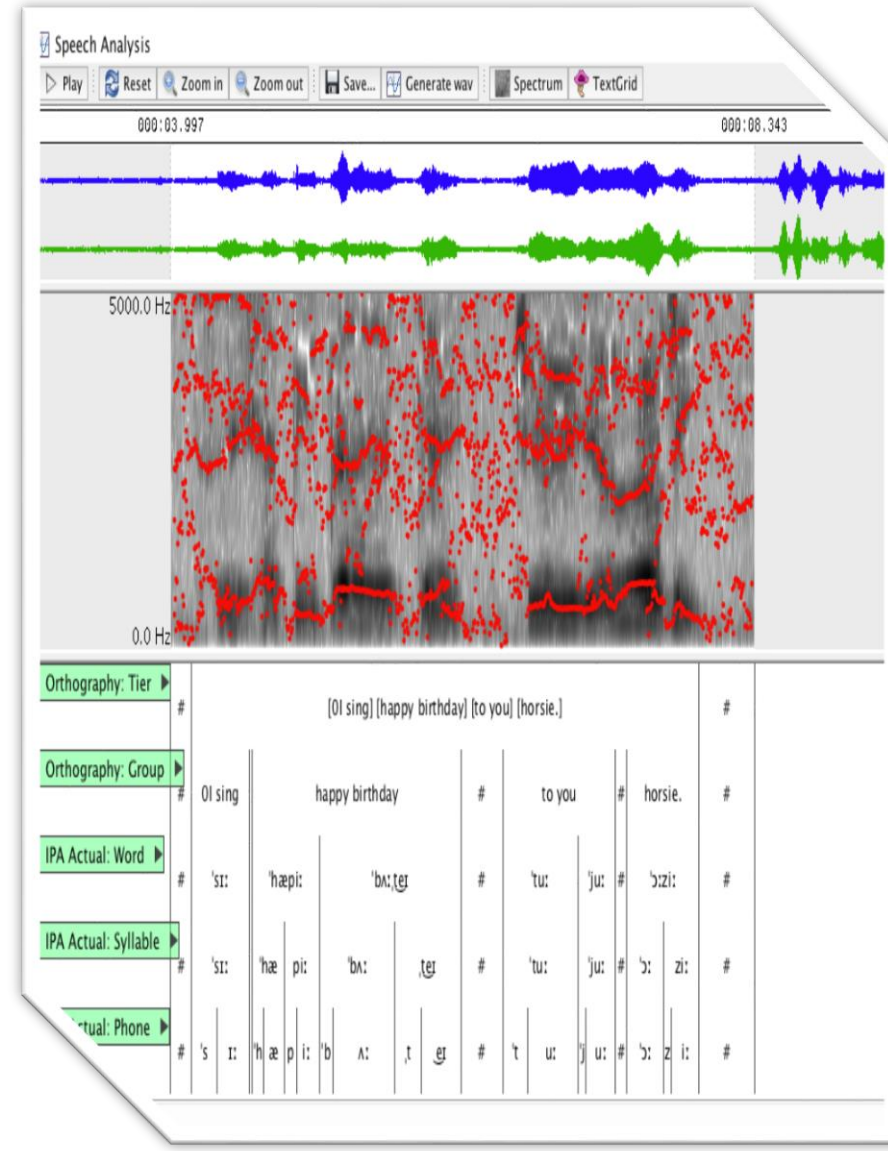# SCAnnAL – An Automatic Speech Corpus Annotator for African Speech Corpora

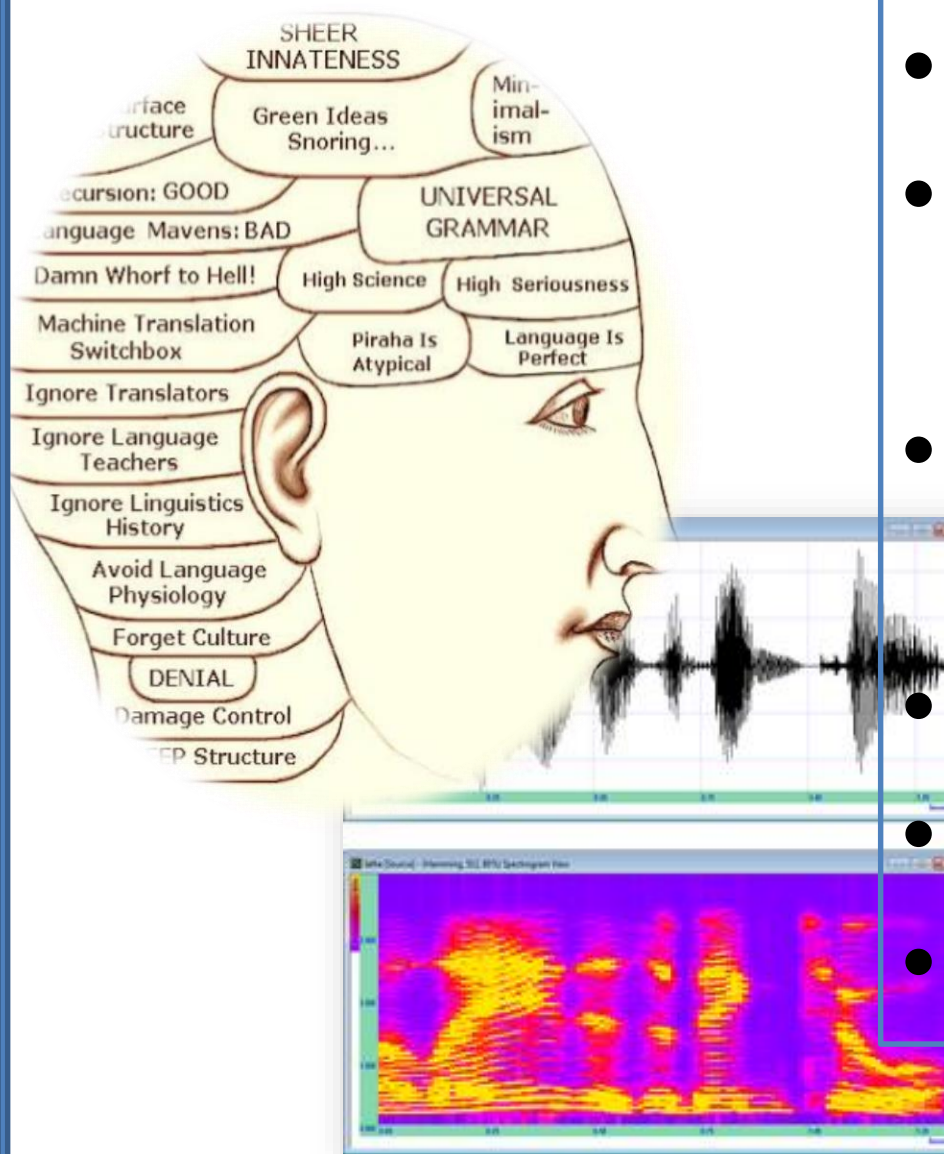## Moses Ekpenyong[1] – Eno-Abasi Urua[2] – Aniefon Akpan[2]

*[1]Department of Computer Science; [2]Department of Linguistics and Nigerian Languages*
(mosesekpenyong@{uniuyo.edu.ng, gmail.com}; enoabasi@gmail.com; daniels.mmefon@gmail.com )

## WHAT IS THE ISSUE?

- Thousands of annotated speech corpora exist worldwide
- The demand for richly annotated corpora is fast growing
- The process accompanying corpora annotation has slowed research progress for African languages
- Current annotation Toolkits do not satisfy the challenges African speech systems present.

## AUDIENCE:

- Linguists
- Fieldwork/language experts
- Computational linguists
- Historians
- Speech Engineers
- Speech Technologists

## OUR GOALS:

- Examine current annotation Toolkits and identify their limitations
- Study the peculiarities of African tone languages
- Automate the annotation process using Signal Processing and NLP
- Adapt automated process to African languages
- Evaluate the annotator for precision

## OUR APPROACH:

- Using Signal Processing, detect the respective speech waveforms
- Segment specified tier(s)
- Accept corresponding transcriptions
- Perform NLP to pre-annotate the transcriptions (e.g., syllabification, ...)
- Render transcriptions (phonemes; words; syllables; sentences; etc.), to respective segments using NLP
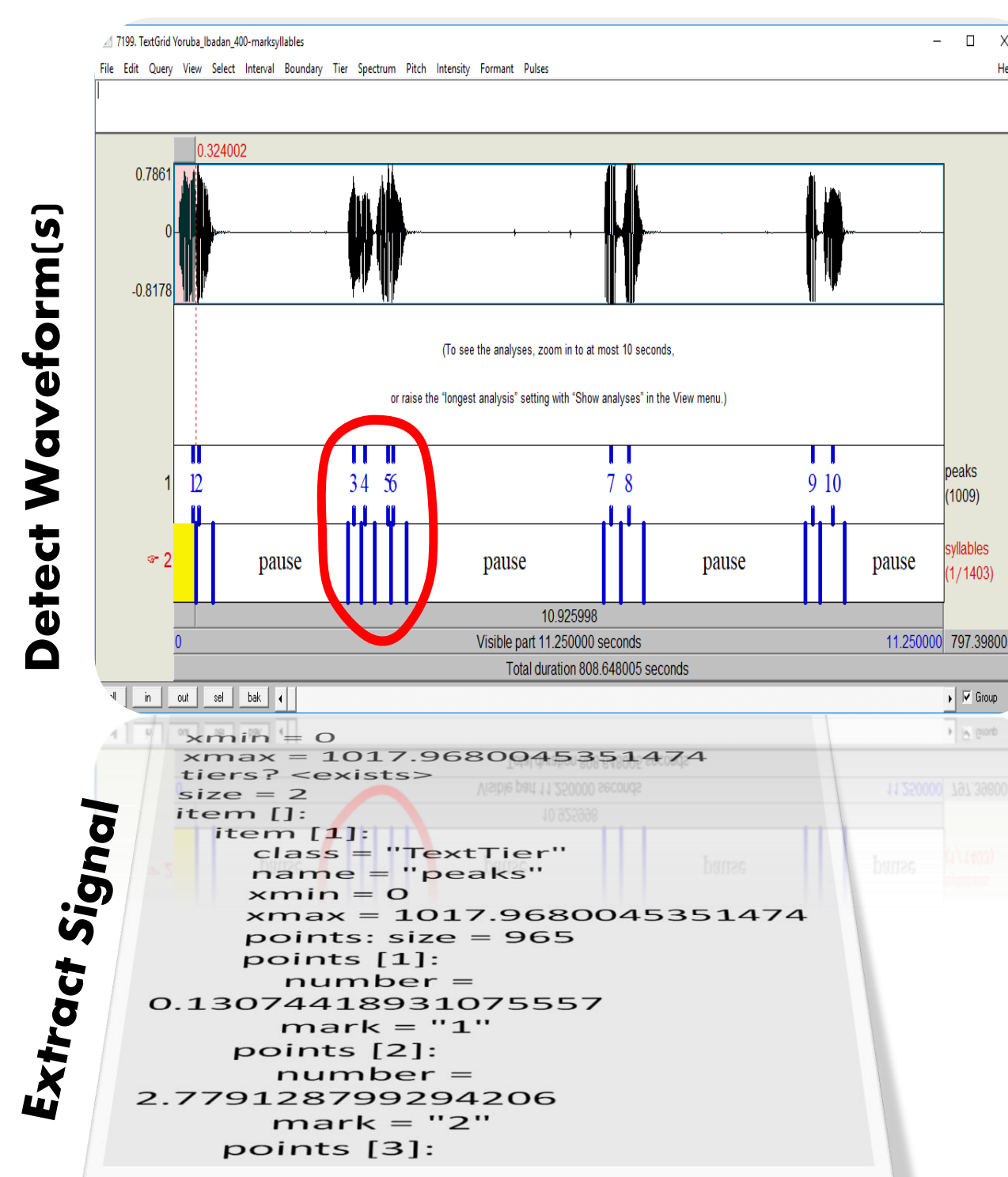- Label segments

## SCAnnAL METHODOLOGICAL WORKFLOW:
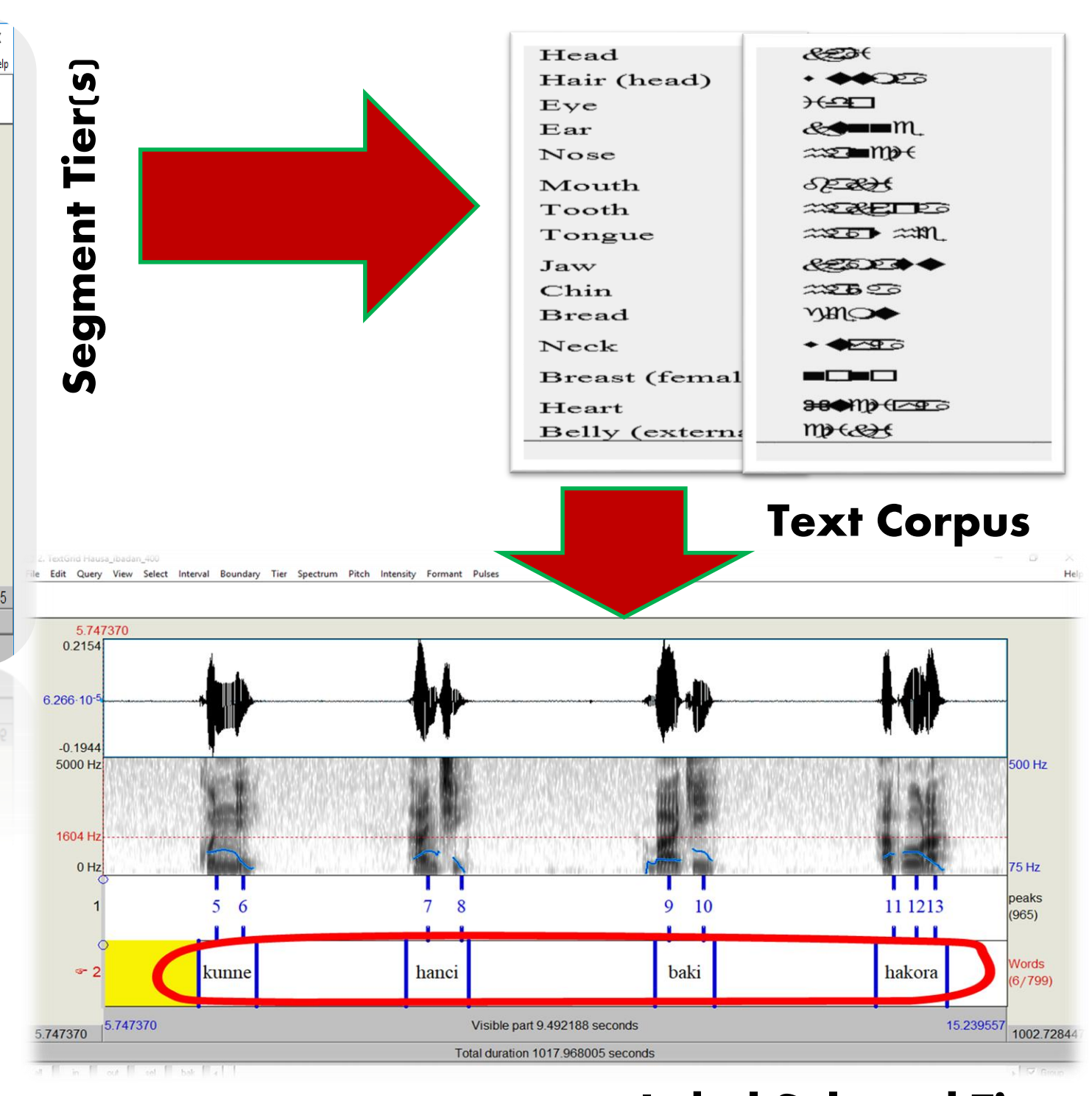


Yoruba Speaker

Hausa Speaker

Detect Waveform(s)

Extract Signal

Segment Tier(s)

Text Corpus

Label Selected Tier