# Interlinearization in ELAN

Han Slöetjes
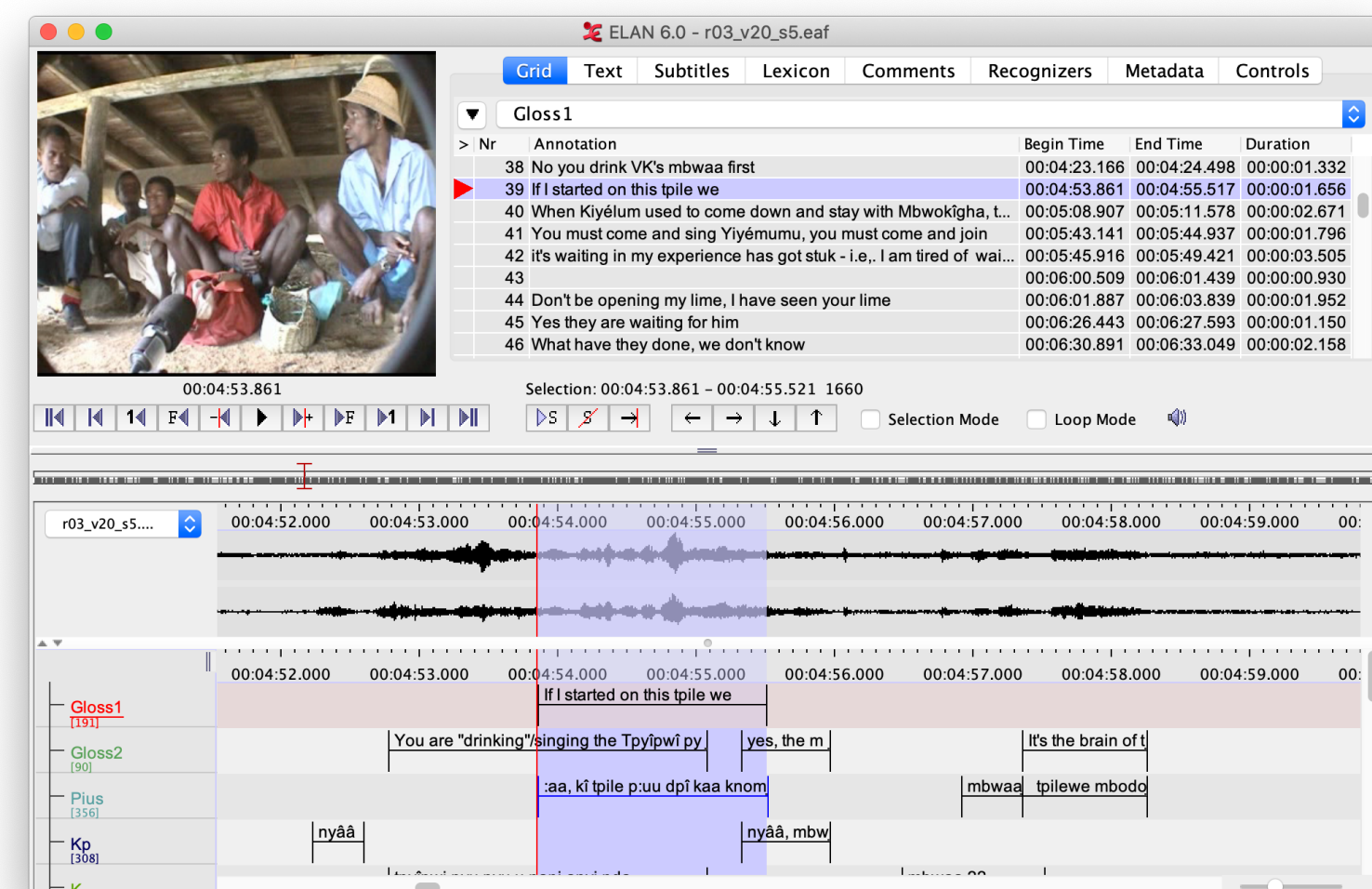
han.sloetjes@mpi.nl
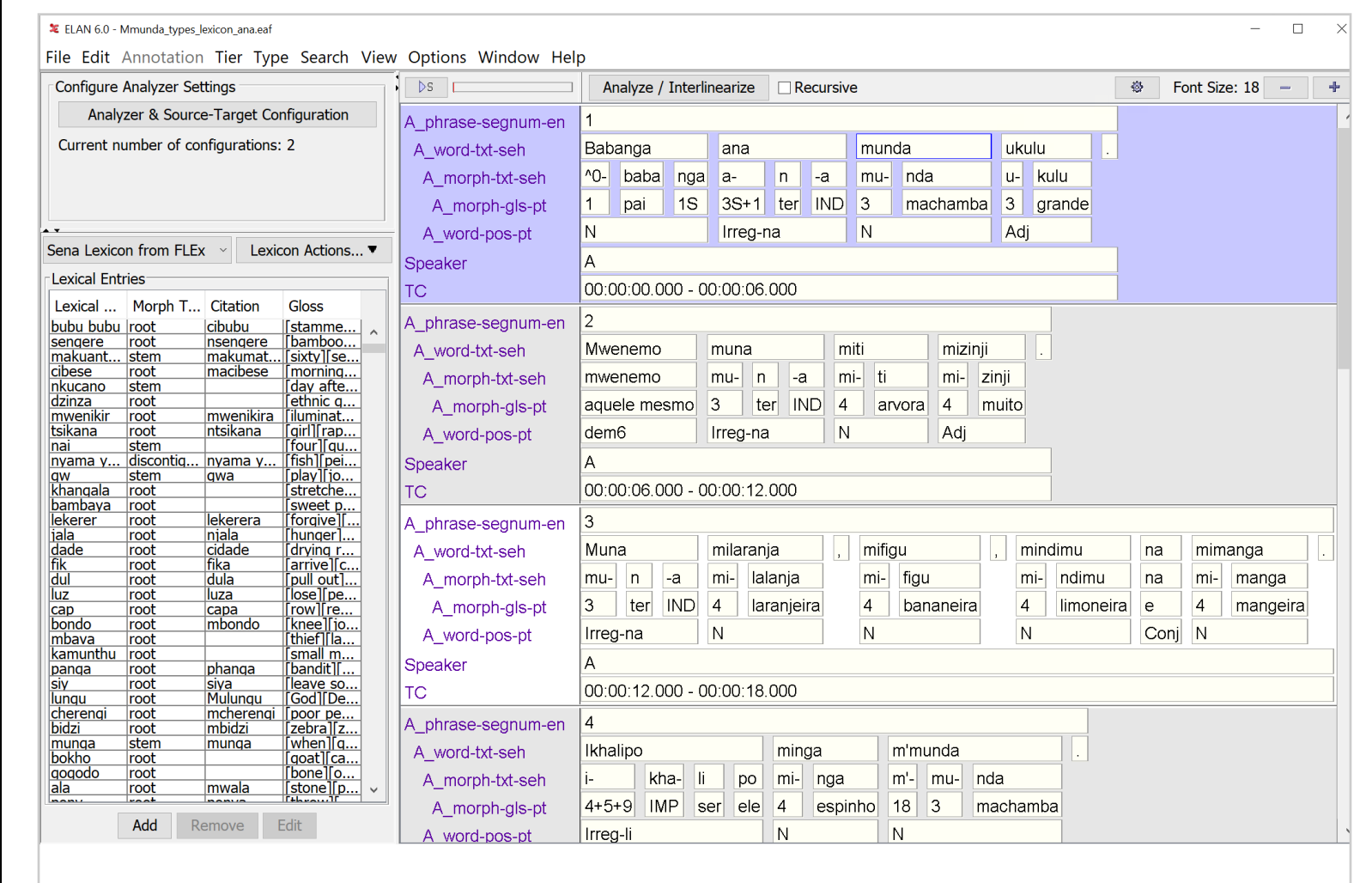
https://archive.mpi.nl/tla/elan

## Introduction

- **ELAN** is a manual annotation tool developed by **The Language Archive/MPI for Psycholinguistics**. It supports multi-tier, multi-speaker, time-linked annotation of audio and video and is applied in many fields of research, language documentation being one of them.



## Background and history

- in language documentation the transcription and translation steps are often followed by morphological parsing and glossing

- it is possible to perform these steps manually in **ELAN**, but the preferred approach is to apply computer-aided methods as offered by **Toolbox**[3] and **FLEx**[4]

- **ELAN** provides import and export functions for the file formats of these tools, so that users can create a toolchain and move their data from one tool to another for the task at hand

- **TLA** once developed the lexicon tool **LEXUS**[1] and prepared to combine it with **ELAN** by way of a new Natural Language Processing module, **LEXAN**[2]

## Interlinearization mode



- new text oriented mode with an *Interlinear Glossed Text (IGT)* style of user interface
- the **LEXAN** modules became extensions of **ELAN** (after end of support for **LEXUS**)
- combined with a new lexicon component
- supports the process of (machine-assisted) parsing and glossing

# Main characteristics

- text oriented (but still linked to the audio!)
- keyboard driven
- optimized for computer assisted morphological parsing and glossing
- text analyzer modules provide '*suggestions*' to be disambiguated by the user
- contains a lexicon editor and viewer
- a flexible system, therefore some configuration is required



Interlinear editor with a play selection button to play the current active phrase.



The suggestions window displays the suggestions in the same layout as the editor. The most frequently selected parses are displayed at the top.

1. Ringersma, J. and Kemps-Snijders, M. (2007). Creating multimedia dictionaries of endangered languages using LEXUS. Interdpeech 2007
2. Stehouwer, H. and Drude, S. (2012). Lexan: A lexical annotation framework for ELAN. Talk presented at LREC 2012

# Text analyzers

- are modules that receive text as input (an annotation) and produce output for a single or for multiple annotations
- can be added as extensions through an API
- can be connected to a lexicon



Linking analyzers to tier types and configuration of an analyzer.

- built-in analyzers:
  - whitespace splitter
  - morphological parser
  - gloss analyzer

- the parser and glosser:
  - require access to a lexicon
  - keep record of choices made to improve suggestions
  - are language independent
  - the parser is implemented as a finite state machine

3. Field Linguist's Toolbox, https://software.sil.org/toolbox/
4. FieldWorks Language Explorer, https://software.sil.org/fieldworks/
5. Lexicon Interchange Format, https://github.com/sillsdev/lift-standard
6. http://corpafroas.tge-adonis.fr

# Lexicon editor



- multiple lexicons can be created and linked
- lexicon entries have a few predefined fields
- custom fields can be added
- sorting of entries according to a custom sort order
- the structure of a lexical entry is similar to the **LIFT**[5] format
- import lexicons from **Toolbox**, **FLEx** (**LIFT**) and **CorpAfroAs**[6] format
- export a lexicon to **LIFT** format