# PDF-to-Text Reanalysis for Linguistic Data Mining

**Michael Wayne Goodman**    **Ryan Georgi**    **Fei Xia**

Linguistics Department
University of Washington
Seattle, WA, USA
{goodmami, rgeorgi, fxia}@uw.edu

## Motivation

- **Many online academic publications are in the form of PDF files, and extracting information from them can be very useful:**
  - Ex: Millions of online linguistic documents that contain language data.

- **Extracting semi-structured text from PDF files is not trivial:**
  - Existing PDF-to-text converters have limitations

- **We develop a system (called Freki) for the reanalysis of PDF-extracted text:**
  - It performs block detection, respacing and tabular data analysis.

## Interlinear Glossed text



## The XY-cut algorithm
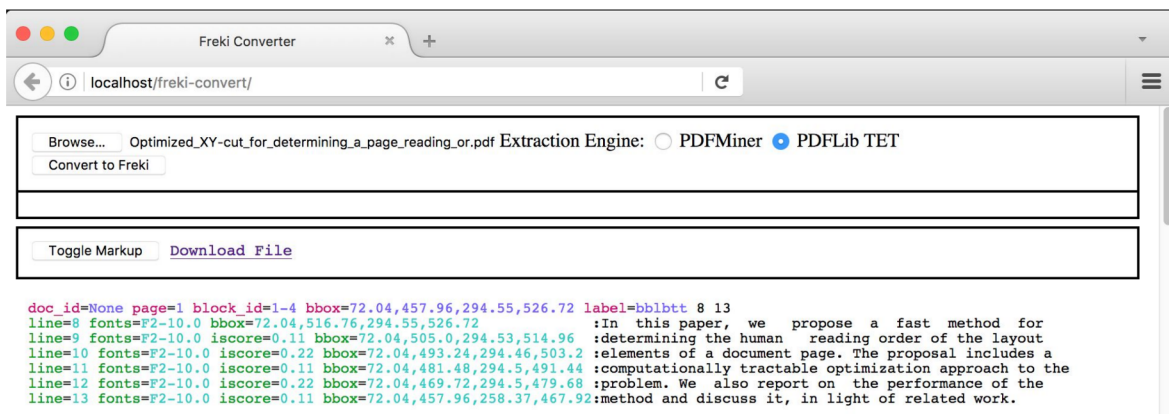


Initial Page    First Cut    Second Cut    Third Cut

## Our approach

- **Built on top of existing PDF-to-text converters that produce XML format.**

- **Added functionality:**
  - Block detection
  - Respacing and tabular data analysis
  - Changes to the XY-cut algorithms

- **Used a new output format for less verbosity and more readability**

## The interface



## Using Freki

- **We ran Freki on millions of documents from the Web which may contain interlinear gloss text (IGT) data.**

- **From the extracted text, we performed IGT detection, language identification, etc.**

- **This leads to an extension to the Online Database of INterlinear text (ODIN), which contains IGT for thousands of languages.**

## Summary

- **PDF-to-text is not trivial, especially if we want to preserve the structure within the text.**

- **Our system, Freki, builds on top of existing converters and focuses on identifying blocks and respacing.**

- **The package is available at** github.com/xigt/freki