# LeXmart
## A Smart Tool for Lexicographers

A. Simões, A. Salgado, R. Costa & J. J. Almeida

IPCA
INSTITUTO POLITÉCNICO
DO CÁVADO E DO AVE

2Ai APPLIED
ARTIFICIAL
INTELLIGENCE
LABORATORY

CLUNL

Universidade do Minho
Escola de Engenharia

# Outline

# Motivation

- Academy of Sciences Portuguese dictionary
- Published in paper format (2001)
- A new goal for 2019–202?: publish an updated version online
- The original dictionary was only available in PDF format
- A conversion to a computer readable format was required
- Need for a tool to allow concurrent edition of the dictionary

# LeXmart Birth

- Develop a process to convert the PDF format into XML (TEI)
- Split the dictionary by entries, in different files
- Import entries into a database that:
  - should be aware of XML structure
  - should be able to query entries using XML structure
  - allow quick editing on XML aware IDEs
- eXist-DB was chosen
  (easy to integrate with Oxygen XML Editor)
- A set of small tools to assess current dictionary status:
  - validate entries
  - edit entries
  - show reports on entry classifications/grammar information

2Ai APPLIED
ARTIFICIAL
INTELLIGENCE
LABORATORY

http://lexmart.eu

# Architecture

- A collection over eXist-DB for the dictionary:
  - Each dictionary entry is an individual XML file
  - Entries are codified using the TEI Dictionary module
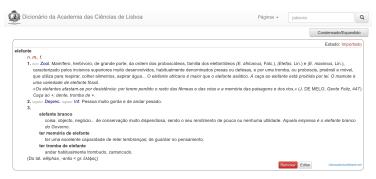  - TEI headers are not being currently stored

# Architecture

- A collection over eXist-DB for the dictionary:
  - Each dictionary entry is an individual XML file
  - Entries are codified using the TEI Dictionary module
  - TEI headers are not being currently stored
- A collection for the application:
  - Development in XQuery 3.0 and XPath 3.0
  - TEI embedded into HTML directly, using CSS

# Architecture

- A collection over eXist-DB for the dictionary:
    - Each dictionary entry is an individual XML file
    - Entries are codified using the TEI Dictionary module
    - TEI headers are not being currently stored
- A collection for the application:
    - Development in XQuery 3.0 and XPath 3.0
    - TEI embedded into HTML directly, using CSS
- A collection for the CMS pages

**2Ai** APPLIED ARTIFICIAL INTELLIGENCE LABORATORY

# End-User Tools

- ▶ The typical use of a dictionary
- ▶ Quite efficient, querying the `orth` TEI elements.



search by "elefante" (*elephant*)

# End-User Tools

- ▶ not possible on paper dictionaries
- ▶ allow querying by the definition
- ▶ queries all entry definition information (slower)



search by "fio ferro" (*iron thread → wire*)
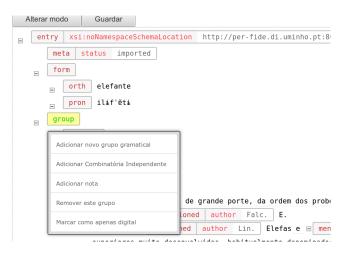
# Lexicographic Support Tools

## Entry editor

- ▶ Based on Xonomy XML editor (JavaScript based)
- ▶ Entries as fetched and saved directly from/into eXist-DB
- ▶ Configurable accordingly with the XML standard being used
- ▶ Contextual menus

# Lexicographic Support Tools

Entry editor

# Lexicographic Support Tools

Entry creation and deletion

► When viewing an entry, lexicographer can delete it (after confirmation)

► Adding a new entry will require the headword and:
  ► guarantee no other entry with the same headword
  ► create a stub XML document, adding it to the database
  ► open the editor with that document

▶ Allow lexicographers to annotate entries and/or senses:

    ▶ imported → directly imported from the paper dictionary
    ▶ new → a new entry, created directly using LeXmart
    ▶ edited → the entry was edited, but not yet considered final
    ▶ revised → the entry can be made available publicly)

    (together with this information, a timestamp is recorded)

▶ Annotate entries as to be used only on digital version

APPLIED
ARTIFICIAL
INTELLIGENCE
LABORATORY

# Lexicographic Support Tools

Filters and Statistics

- Entries have diverse type of metainformation:
  - grammatical information
  - geographic source
  - register type
  - domain knowledge

# Lexicographic Support Tools

- ► Entries have diverse type of metainformation:
  - ► grammatical information
  - ► geographic source
  - ► register type
  - ► domain knowledge
- ► Filtering allows reports for specific kinds:
  - ► a list of terms from a specific geographic location
  - ► terms from a specific knowledge area
  - ► etc.

- ▶ Entries have diverse type of metainformation:
    - ▶ grammatical information
    - ▶ geographic source
    - ▶ register type
    - ▶ domain knowledge
- ▶ Filtering allows reports for specific kinds:
    - ▶ a list of terms from a specific geographic location
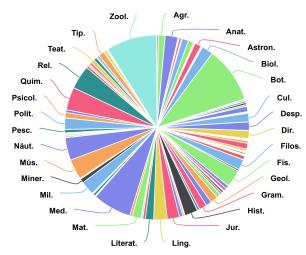    - ▶ terms from a specific knowledge area
    - ▶ etc.
- ▶ Statistics:
    - ▶ allow to evaluate relevance of classes
      *How many entries in the cutlery domain?*
    - ▶ Allow to validate possible duplicated tags

# Lexicographic Support Tools

Filters and Statistics

▶ Allow to create custom reports on work evolution:
  - ▶ list of entries accordingly with their state
    (imported, new, edited, and revised)
  - ▶ digital only entries
  - ▶ entries from a specific defined list
    (given a text file with the list of terms)

▶ Validation:

  ▶ LeXmart validates for well-formedness of the XML
  ▶ LeXmart does not validate for DTD/Schema compliance in run time
  ▶ Reports can be requested on Schema compliance

- Validation:
  - LeXmart validates for well-formedness of the XML
  - LeXmart does not validate for DTD/Schema compliance in run time
  - Reports can be requested on Schema compliance
- Backups:
  - eXist-DB backups are dumps of XML files
  - These files can be stored in a control version software
  - Using GIT allows us to have backups and version control easily

# Content Management System

- Dictionaries are not just a list of definitions:
  - introduction
  - explanation of micro and macro structure
  - other decisions...
- Information to be shared by lexicographers:
  - rules when editing entries
  - general guidelines

# Content Management System

▶ A simple CMS system was developed:

  ▶ stored as independent collection on eXistDB
  ▶ edited on a WYSIWYG editor (TinyMCE)
  ▶ stored as XHTML
  ▶ easy to display directly on a browser

# RESTful API

- ▶ eXist-DB provides a RESTful API to:
  - ▶ fetch/store entries directly
  - ▶ query collections using XQuery
- ▶ On top of this, a generic user-centric API will be developed:
  - ▶ simple access to search/reverse search
  - ▶ filtering accordingly with entries metadata

APPLIED
ARTIFICIAL
INTELLIGENCE
LABORATORY

# Conclusions

▶ Solution based on Web-technologies

▶ Based on XML-aware database

▶ Extensible

▶ Available freely on GitLab (see `http://lexmart.eu`)

▶ Being used on a real dictionary project

# Future Work

- Translate the software
- Add configuration properties
- Stick to TEI and/or TEI Lex-0 standards
- Create "print-friendly" formats (XSL-FO or LaTeX)
- Integrate with other technologies:
  - ontologies and triple stores
  - wordnets

# LeXmart

## A Smart Tool for Lexicographers

A. Simões, A. Salgado, R. Costa & J. J. Almeida