

# Interlinking Lexicographic Data in the MORDigital Project

*Anas Fahad Khan*, CNR-ILC, *Ana Salgado*, Centro de Linguística da Universidade Nova de Lisboa (CLUNL) & Academia das Ciências de Lisboa, *Margarida Ramos*, CLUNL, *Sara Carvalho*, CLUNL & Centro de Línguas, Literaturas e Culturas, *Laurent Romary*, Automatic Language Modelling and ANALysis & Computational Humanities Inria de Paris (ALMANACH), *Bruno Almeida*, CLUNL, *Mohamed Khemakhem*, ArcaScience, *Raquel Silva*, CLUNL, *Toma Tasovac*, Belgrade Center for Digital Humanities (BCDH), *Rute Costa*, CLUNL.

# Introduction

---

- In this talk we will introduce the Portuguese national project **MORDigital** and present some updates on its current progress.
- MORDigital aims to make a historic Portuguese-language dictionary, *Dicionário da Língua Portuguesa* aka *Morais*, available as a digital resource.
- The project brings together some of the latest innovations in **computational lexicography**, the **digital humanities** and **linguistic linked data** including work on modelling lexicographic resources using standards such as **OntoLex**, **TEI Lex-0** and **LMF**.
- It also innovative in its creation and use of pipelines for converting retrodigitised dictionaries into computational lexical resources.

# The *Diccionario da Lingua Portuguesa*

---

- As the first Portuguese monolingual dictionary *Morais* was instrumental in normalising the language and become the model for subsequent Portuguese language dictionaries.
- It was influenced by other modern language dictionaries published in Europe in the 16th and 17th centuries during the age of Enlightenment.
- Authorship of the 1st edition is attributed to Rafael Bluteau, a Portuguese priest and lexicographer, whose Portuguese-Latin Vocabulary (10 vols., 1712-1728) was the basis for *Morais*.
- *Morais* directly oversaw the 2nd (1813) and 3rd (1823) editions, which greatly overhauled the dictionary.

# The *Diccionario da Lingua Portuguesa*

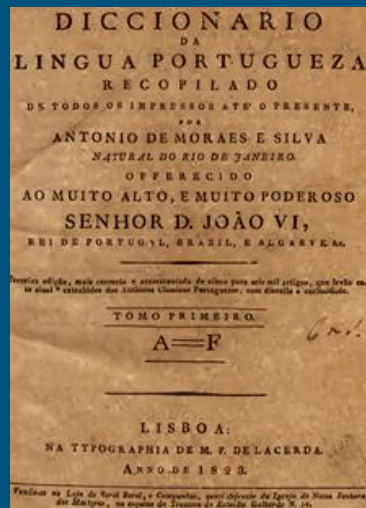
Frontispiece of *Morais* dictionaries (1789, 1813, 1823)



Two volumes  
A to K, 752 pp.  
L to Z, 541 pp.



Two volumes  
A to E, 889 pp.  
F to Z, 886 pp.

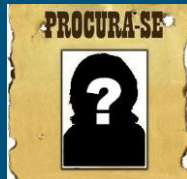


Two volumes  
A to K, 952 pp.  
L to Z, 872 pp.

# Who was Moraes?

---

- The dictionary in question is usually referred to by the name of its compiler, the renowned Brazilian lexicographer **ANTÔNIO DE MORAIS SILVA** (1757?–1824)
- Born in Rio de Janeiro, Moraes graduated in Civil and Canon Law from the University of Coimbra in Portugal.
- After being accused of heresy by the Inquisition, he fled to England and devoted himself to the study of languages; it was here that he planned the structure of his future dictionary
- Moraes moved to Brazil in 1794 where he entered the judiciary and held the position of judge in the Bahia Court of Appeal (a position from which he soon resigned). He subsequently moved to Muribeca, in Pernambuco, where he lived until his death on April 11, 1824
- We can't find any images of Moraes himself!



# MORDigital

---

- **MORDigital – Digitalização do Dicionário da Língua Portuguesa de António de Moraes Silva** [PTDC/LLT-LIN/6841/2020] is a project financed by the Portuguese National Funding agency through the FCT – Fundação para a Ciência e Tecnologia
- Although it is a Portuguese national project, it also includes collaborators from all over Europe.
- The main aim of the project is to convert **three editions** (1789; 1813; 1823) of *Moraes* into a structured lexical resource in both TEI-XML (based on the ISO LMF standard) and in RDF (based on the OntoLex-Lemon model and its recent extensions).
- These editions will also be made available via an online interface on the website (at the moment only PDFs are available).

# MORDigital - The website

The screenshot displays the MorDigital website. On the left, a large text box contains the following information:

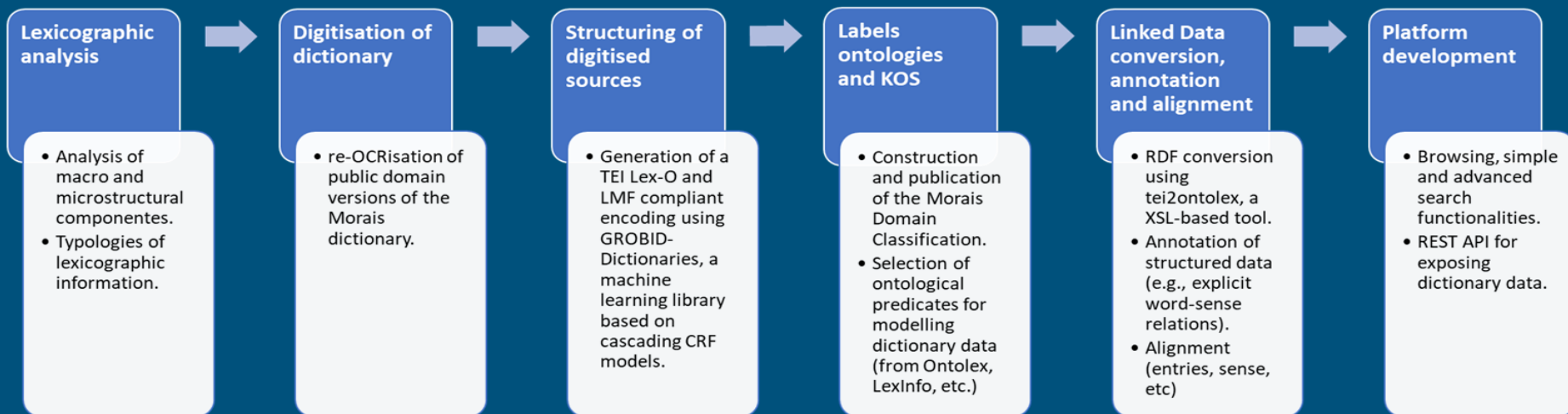
**MORDigital** is a lexicographic project that aims to produce digital versions of the first three editions (1789; 1813; 1823) of the *Diccionario da Lingua Portuguesa* by António de Morais Silva.

Below this text, the URL <https://mordigital.fcsh.unl.pt/> is provided.

The website header includes the MorDigital logo and navigation links: ABOUT, DICTIONARIES, RESULTS, EVENTS, NEWS, and TEAM. The main content area features a section titled "About the Project" and a "Team" section with portraits of the project members. On the right, a preview of the "Dictionary - 1st Edition (1789)" is shown, displaying the title page and several pages of the dictionary's content.



# Workflow





- An established tool with user-friendly interface.
- It creates editable, searchable files. Produces several output formats (RTF, DOCX, PDF, HTML, XML, etc.).
- The tool preserves typographical features.



Noise produced by the OCRization	Description	Example
Archaic characters	The <b>long /s/</b> is recognised as a lower-cased <i>/f/</i> , or <i>/j/</i> if the original is typed in italic	<b>f</b> instead of <b>f</b> <b>j</b> instead of <b>f</b>
Case-sensitivity	Lower- and upper-case characters can be mixed up	<b>I</b> instead of <b>i</b>
Wrong characters inserted	Characters are misrecognised as wrong characters	<b>rn</b> instead of <b>m</b>
Segmentation errors	Different spacings between lines, words. Misrecognition of white-spaces	<b>temtres</b> instead of <b>tem tres</b>
Ligature	A combination of two or three characters set as a single glyph	<b>fi</b>
Ink spots	Text has ink spots, which prevents both human and machine from reading	

# Structuring of Digital Editions

---

- For the structuring of the digital editions we are using **GROBID-Dictionaries**
- This is a machine learning library for structuring digitised lexical resources and entry-based documents with encyclopedic or bibliographic content.
- In particular, it allows the automatic parsing, extraction and structuring of lexical information from PDFs.
- GROBID-Dictionaries takes as input lexical resources digitised in PDF format and generates a TEI-encoded hierarchy of the different text structures which it has recognised.

# TEI and OntoLex

---

- The TEI-XML sources will subsequently be converted to OntoLex (both the original model and its follow-up modules) using an XSLT stylesheet.
- Having an RDF version available allows us to make the dictionary available using both a dedicated platform and via a SPARQL endpoint.
- In addition, the RDF versions of each edition of the dictionary will be published on the LLOD cloud. This will be an important contribution to adding Portuguese language resources to the cloud.
- At the end of the project, based on our experiences, we will draft technical guidelines to help lexicographers and digital humanists carry out these tasks.

# An Example Entry

META'STASE, ou *Metastasis*, f. f. *Med.* de-  
 geração de huma doença em outra, especie de  
 Crife. § *na Rhet.* figura pela qual o Orador at-  
 tribue alguma coisa a outrem, desonerando-se  
 della.

```
<entry xmlns="http://www.tei-c.org/ns/1.0"
type="monolexicalUnit" xml:lang="pt"
xml:id="MORAIS_1.metastase">
  <form type="lemma">
    <orth>METÁSTASE</orth>
    <pc>, ou</pc>
    <form type="variant">
      <orth>Metastasis</orth>
    </form>
    <pc>,</pc>
    <gramGrp>
      <gram type="pos">
norm="NOUN">s.</gram>
      <gram type="gen">f.</gram>
    </gramGrp>
    <sense xml:id="MORAIS_1.metastase_1">
      <usg type="domain">Med.</usg>
      <def>degeneração de huma doença em
outra, espécie de Crise</def>
    </sense>
    <sense xml:id="MORAIS_1.metastase_2">
      <pc>na</pc>
      <usg type="domain">Rhet.</usg>
      <def>figura pela qual o Orador attribue
alguma coisa a outrem , desonerando-se
della.</def>
    </sense>
  </entry>
```

# Organising knowledge

Morais	METALABEL	Encyclopedia of Diderot and d'Alembert	Dewey Decimal Classification (DDC)	EuroVoc	BabelNet
Arithm. Arithmetico	arithmetics	Reason Philosophy Science of Nature Mathematics Pure Arithmetics	500 Natural sciences & mathematics 510 Mathematics 513 Arithmetic	BT1 pure mathematics BT2 mathematics BT3 natural sciences	IS A pure mathematics area of mathematics
Geometr. Geometrico	geometry	Reason Philosophy Science of Nature Mathematics Pure Geometry	500 Natural sciences & mathematics 510 Mathematics 516 Geometry	BT1 pure mathematics BT2 mathematics BT3 natural sciences	IS A pure mathematics area of mathematics mathematics
Mathem. Mathematico	mathematics	Reason Philosophy Science of Nature Mathematics	500 Natural sciences & mathematics 510 Mathematics	BT1 natural sciences NT1 applied mathematics NT1 pure mathematics	IS A Science major Universal language Academic discipline formal science

# Organising knowledge

---

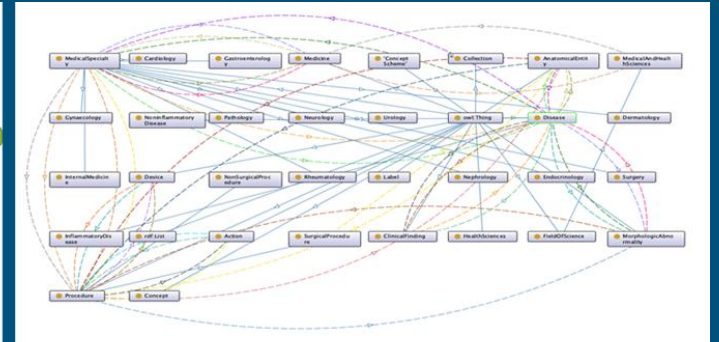
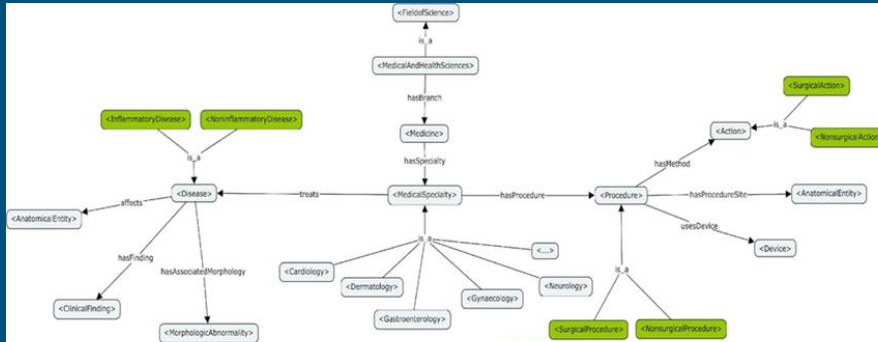
ARITHMETICA, f. f. arte de calcular por algarifmos.

GEOMETRIA, f. f. parte da Mathematica que ensina a conhecer a grandeza, razões, e proporções das grandezas continuas, ou sejam linhas, ou figuras, ou sólidos, ou superficies.

MATHEMATICA, f. f. a sciencia, que ensina a conhecer as grandezas de toda sorte, suas razões, relações, e proporções: *Mathematica mista* (oppõe-se ás puras) a que ensina a applicar os principios de calculo, e geometria aos corpos.



# Organising knowledge



Active ontology: **Protocols** > Individuals by class

Classes (Object properties) Data properties Annotation properties Datatypes Individuals

**Individuals** **NonSurgicalProcedure**

- not Thing
- not Action
- not AnatomicalEntity
- not ClinicalFinding
- not Collection
- not Concept
- not Connector Schema
- not Device
- not Disease
- not FluoroscopicSource
- not ImagingScience
- not Lab
- not MedicalSpecialty
- not Medicine
- not MorphologicalAbnormality
- not Procedure
- NonSurgicalProcedure**
- not TemporalMeasurement
- not Unit

**Annotations: SurgicalProcedure**

Annotations (Usage)

Annotation	Value
rdfl:label	Temporally
procedimento:qnameproc	
rdfl:label	Temporally
procedimento:liturgic	
rdfl:label	Temporally
procedimento:qnameproc	

**Data classes: SurgicalProcedure**

Instances (Usage)

Instance	Value
hasMethod	SurgicalAction
Procedure	

Domain class position

Domain class	Value
hasDevice	Device
hasProcedureWrite	AnatomicalEntity
hasMethod	Action
ProcedureOf	MedicalSpecialty

Instances

Target for key

Domain term

**NonSurgicalProcedure**

Domain class: CP

# Ongoing Work

---

- Start OCR corrections on the two other editions of the dictionary
- Run the GROBID tool iteratively on the corrected output of the OCR of the first edition to ensure a correct TEI-XML (LMF) encoding of the different components of single entries (e.g., authoritative examples, collocations, )
- Start testing the XSLT transformation to OntoLex and its extensions
- In addition throughout the project our approach foresees a rigorous linguistic treatment of the source texts. This will make it possible to organise and structure the lexicographic components, and to elicit lexical relationships between various elements.
- We also propose combining semasiological and onomasiological approaches in our treatment of the different editions of Moraes. For this we foresee the inclusion of ontologies (e.g. diasystematic marking, namely domain labels, registers and part of speech categories).

# Thank You! Obrigado! Ačiū!

---

Anas Fahad Khan: [fahad.khan@ilc.cnr.it](mailto:fahad.khan@ilc.cnr.it)

Ana Salgado: [anasalgado@fcsh.unl.pt](mailto:anasalgado@fcsh.unl.pt)

Rute Costa: [rute.costa@fcsh.unl.pt](mailto:rute.costa@fcsh.unl.pt)

Sara Carvalho: [sara.carvalho@ua.pt](mailto:sara.carvalho@ua.pt)

Laurent Romary: [laurent.romary@inria.fr](mailto:laurent.romary@inria.fr)

Bruno Almeida: [brunoalmeida@fcsh.unl.pt](mailto:brunoalmeida@fcsh.unl.pt)

Margarida Ramos: [mvramos@fcsh.unl.pt](mailto:mvramos@fcsh.unl.pt)

Mohamed Khemakhem: [medkhemakhemfsegs@gmail.com](mailto:medkhemakhemfsegs@gmail.com)

Toma Tasovac: [ttasovac@humanistika.org](mailto:ttasovac@humanistika.org)

Raquel Silva: [raq.silva@fcsh.unl.pt](mailto:raq.silva@fcsh.unl.pt)

Toma Tasovac: [ttasovac@humanistika.org](mailto:ttasovac@humanistika.org)