

Introduction to CLTK (Classical Language ToolKit)  
Eleftheria Chatziargyriou & Clément Besnier  
06/11/2019



## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

# **CLTK: PHILOSOPHY AND ORGANIZATION**

## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

- **Free and Open-Source Python** library
- Founded in 2014 by **Kyle P. Johnson**
- Provides **NLP<sup>1</sup> tools** for **historical languages**
- **Shares** a **high-quality code** for **academic research**

Co-maintainers: Patrick J. Burns and Kyle P. Johnson.

---

<sup>1</sup>NLP: Natural Language Processing

Main goals:

1. Compile analysis-friendly corpora
2. Collect and generate linguistic data<sup>2</sup>
3. Act as a free and open platform for generating scientific research

---

<sup>2</sup>[https://github.com/cltk/latin\\_models\\_cltk](https://github.com/cltk/latin_models_cltk)

## 1 CLTK: philosophy and organization

- Overview
- **NLP Tools**
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

# CLTK AMONG OTHER NLP TOOLS IN PYTHON

- NLP: SpaCy<sup>3</sup>, NLTK<sup>4</sup>, StanfordNLP<sup>5</sup>
- Python is a programming language widely used by researchers

---

<sup>3</sup><https://spacy.io/>

<sup>4</sup><https://www.nltk.org/>

<sup>5</sup><https://stanfordnlp.github.io/stanfordnlp/>



## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- **Historical Languages**
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

- Early Antiquity: Sumerian, Akkadian, Old Egyptian, etc
- Late Antiquity: Ancient Greek, Latin, Sanskrit, Classical Chinese, Gothic, etc
- Middle Ages: Medieval Latin, Coranic Arabic, Koine, Old and Middle High German, Old Norse, etc

- Handles languages written and spoken before Gutenberg
- Documents written in these languages have specific features:
  - ▶ often relatively small and fragmentary surviving texts
  - ▶ spelling not normalized
  - ▶ diachronic component of languages must be handled
  - ▶ no more living speakers
  - ▶ no more produced texts
- Expert skills needed

## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

- Decentralization
- Disintermediation
- Extensibility
- Standardization
- Simplicity

- Transparency
- Inclusion
- Multi-disciplinary
- Mutual benefit

- Free
- MIT license<sup>6</sup>, you can share and reuse it, even for commercial code
- Inclusion
- Multi-disciplinary

---

<sup>6</sup><https://choosealicense.com/licenses/mit/>

How to cite the project:

```
@Misc{johnson2014,  
  author = {Kyle P. Johnson et al.},  
  title = {CLTK: The Classical Language Toolkit},  
  howpublished = {\url{https://github.com/cltk/cltk}},  
  note = {{DOI} 10.5281/zenodo.<current_release_id>},  
  year = {2014--2019},  
}
```

You can also cite the precise contributors<sup>7</sup> if you use a specific module.

To this day, CLTK has been cited more than 50 times<sup>8</sup>

---

<sup>7</sup><https://github.com/cltk/cltk/blob/master/contributors.md>

<sup>8</sup>from Google scholar



# **CLTK - CODE AND CONTRIBUTION**

## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

# WHAT CAN CLTK DO?

- Corpora importing
- Text preprocessing
  - ▶ File Parsing
  - ▶ Orthographic Normalization
  - ▶ ASCII/Unicode Conversion
  - ▶ Stopword Filtering
- Text processing
  - ▶ Syllabification
  - ▶ Syllable/Word Stressing
  - ▶ Phonetic Indexing
  - ▶ Word/line Tokenization
  - ▶ IPA Transcription
  - ▶ Lemmatization
  - ▶ Stemming
  - ▶ POS Tagging
  - ▶ Poetry Scansion
  - ▶ Named Entity Recognition

The **basic tools** to make **analysis** and **automatic tasks** on a language like text summarisation, question answering, machine translation, etc.

# CURRENTLY SUPPORTED LANGUAGES

Corpora: Bengali, Chinese, Coptic, Egyptian, Gujarati, Hebrew, Javanese, Malayalam, Odia, Old Church Slavonic, Old Swedish, Pali, Persian, Prakrit, Telugu, Tibetan, Urdu

		<b>Akkadian</b>	Arabic	Hindi	<b>Greek</b>	<b>Latin</b>	Marathi	Middle English	<b>Middle High German</b>	Middle Low German	Old English	Old French	Old Norse	Punjabi	Sanskrit
Corpora	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
Stoplist	•	•	•	•	•	•	•	•	•		•	•	•	•	•
Sentence tokenizer					•	•						•			
Word tokenizer	•	•			•	•		•	•			•	•		•
Stemmer	•					•		•	•			•			•
Lemmatizer					•	•		•				•			
POS tagger					•	•			•	•	•		•		
Prosody tagger					•	•			•				•		
NER					•	•						•			

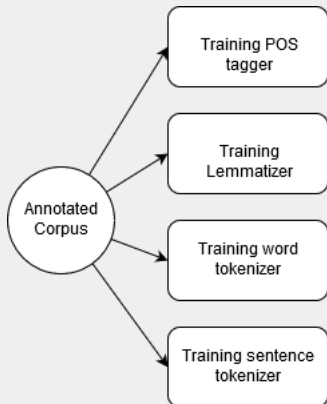
## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

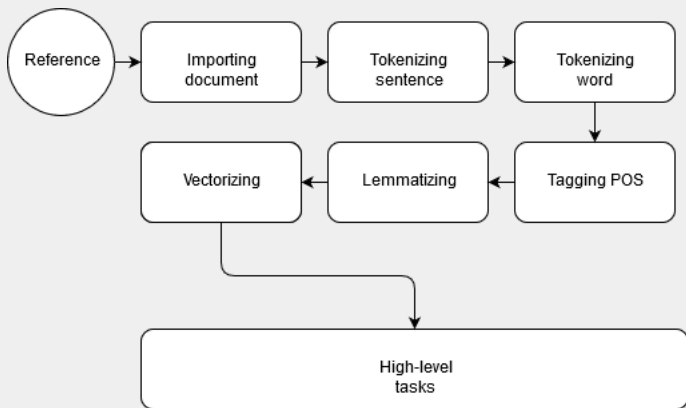
- Code
- **Pipelines**
- Contribution

Building a Text Analysis Pipeline for Classical Languages (Burns 2019, p. 33)





# CLTK PIPELINES



## 1 CLTK: philosophy and organization

- Overview
- NLP Tools
- Historical Languages
- High Quality Code

## 2 CLTK - Code and Contribution

- Code
- Pipelines
- Contribution

- Collaborative effort / open to a virtually infinite talent pool
- Avoid “re-inventing the wheel”
- Closer to the needs of the community
- Constant patches
  - ▶ Bugs are quickly resolved
  - ▶ New features are constantly developed
- Transparency of development
- Generally results in safer software
- Easily customizable

# WHY CONTRIBUTE?

- Best reason: ensure scientific reproducibility of your research
- Expand your skill set
- Give back to the community
- Open Source culture
- It's Fun!

# HOW TO CONTRIBUTE

- You can check out the **CLTK tutorials**<sup>9</sup> and **docs**<sup>10</sup>
- Take a look at the **open issues**<sup>11</sup> or simply make **your own contribution**<sup>12</sup>.
- Don't hesitate to ask for help in the **IRC channel**<sup>13</sup>!

---

<sup>9</sup><https://github.com/cltk/tutorials>

<sup>10</sup><http://docs.cltk.org>

<sup>11</sup><https://github.com/cltk/cltk/issues>

<sup>12</sup><https://github.com/cltk/cltk/pulls>

<sup>13</sup><https://gitter.im/cltk/cltk>

- Open for all university students
- You can work on an open source project for the summer
- CLTK participated for 3 years (2016, 2017, 2018)

## PREVIOUS GSOC PROJECTS

- Additional support Akkadian, Germanic languages, Old and Middle French
- Greek/Latin Backoff lemmatizers.
- Support of more synonyms, translations and word embeddings for Greek and Latin
- Annotation support for the CLTK Archive

- Hosted by Github <https://github.com/cltk/cltk>
- 2723 commits
- 83 contributors
- 73 watchers, 560 stars, 283 forks
- 45 releases, current version 0.1.112
- Code coverage 89



- Digital tools can be used to aid academics and speed up mundane and well-defined processes
- Classical languages have their own unique set of challenges compared to modern languages
- CLTK offers an easy to use and well-documented API for Classical Natural Language Processing




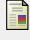
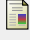

Thank you for your attention!<sup>14</sup>

---

<sup>14</sup>You can now get stickers!

- Krauwer, S. 2003. "The Basic Language Resource Kit (BLARK) as the First Milestone for the Language Resources Roadmap." Proceedings of the 2003 International Workshop on Speech and Computer (SPECOM 2003) : 8-15; cf. also Passarotti in SALT MIL 2010: 29.
- Building a Text Analysis Pipeline for Classical Languages, <https://www.degruyter.com/view/books/9783110599572/978311059010/9783110599572-010.xml>

- Eleftheria Chatziargyriou, email address: [ele.hatzy@gmail.com](mailto:ele.hatzy@gmail.com), Github: [Sedictious](#)
- Clément Besnier, email address: [clemsciences@aol.com](mailto:clemsciences@aol.com), personal website: [clementbesnier.fr](http://clementbesnier.fr), Twitter: [clemsciences](#), Github: [clemsciences](#)

-  PATRICK J BURNS, *CONSTRUCTING STOPLISTS FOR HISTORICAL LANGUAGES*, DIGITAL CLASSICS ONLINE (2018), 4–20.
-  THIBAUT CLÉRICE, *CAPTAINS TOOLKIT, DIGITAL EDITING AND DATA REUSE*, MEDIEVALES-PARIS- **73** (2017), NO. 73, 115–131.
-  OKSANA DEREZA, *LEMMATISATION FOR UNDER-RESOURCED LANGUAGES WITH SEQUENCE-TO-SEQUENCE LEARNING: A CASE OF EARLY IRISH*, EPIC SERIES IN LANGUAGE AND LINGUISTICS **4** (2019), 113–124.
-  LEOPOLD HESS AND CORIEN BARY, *NARRATOR LANGUAGE AND CHARACTER LANGUAGE IN THUCYDIDES: A QUANTITATIVE STUDY OF NARRATIVE PERSPECTIVE*, DIGITAL SCHOLARSHIP IN THE HUMANITIES (2019), FQZ026.
-  TOM KEELINE AND TYLER KIRBY, *AUCEPS SYLLABARUM: A DIGITAL ANALYSIS OF LATIN PROSE RHYTHM*, JOURNAL OF ROMAN STUDIES **109** (2019), 161–204.
-  MARIA MORITZ AND MARCO BÜCHLER, *AN AUTOMATED APPROACH TO MODEL THE TRANSFORMATION PROCESS OF THE REUSE OF BERNARD DE CLAIRVAUX: HOW DO LEXICAL RESOURCES HELP?.*, DH, 2017.