# COLLATINUS : LEMMATIZER AND MORPHOLOGICAL ANALYZER FOR LATIN TEXTS.

Yves Ouvrard
Philippe Verkerk

Digital Classics III: Re-thinking Text Analysis

May 12th 2017

# Who are we ?

Yves Ouvrard
Professor of latin
(retired)

Philippe Verkerk
Physicist
Lab. PhLAM

# What was initially Collatinus ?

- Prepare the wordlist associated to a text
- Reading support of latin texts

# Collatinus 11

- Free and Open program (GNU GPL)
  - written in C++ (Qt 5)
  - resources as plain text files


- Stand-alone software (Mac OS X, Windows, Linux)

  http://outils.biblissima.fr/en/collatinus/

- On-line version (Collatinus 10.2)

  http://outils.biblissima.fr/en/collatinus-web/

# Main functions

- ## Lemmatization
  Analysis
  Reading support

- ## Dictionaries
  - Gaffiot 2016
  - Lewis & Short
  - Georges...

- ## Inflection

| sing. | M. | F. | N. |
|-------|------|-------|-------|
| nom. | dŏctŭs | dōctă | dōctŭm |
| uoc. | dōctĕ | dōctă | dōctŭm |
| acc. | dōctŭm | dōctăm | dōctŭm |
| gen. | dōctī | dōctāe | dōctī |
| dat. | dōctō | dōctāe | dōctō |
| abl. | dōctō | dōctă | dōctō |

- ## Scansion

Ārmă (Ārmā) vĭrūmquĕ cānō (cănō), Trōjāe quī prīmŭs ăb ōrīs (ōrĭs)
-- u-u -- -- - -u u --
Ītălĭām, fātō prŏfŭgŭs, Lāvīnĭăquĕ (Lāvīnĭăquĕ) vēnĭt (vĕnĭt)
-uu- -- uu- --u-u -u
lītŏră, mūlt[um] īll[e] ēt tērrīs jăctātŭs (jăctātŭs) ĕt āltō
-uu -` -` - - -- --- u --
vī sŭpĕrūm sāevāe mĕmŏrēm Jūnōnĭs ŏb īrăm;
- uu- -- uu- --u u -u
mūltă (mūltā) quŏqu[e] (quŏqu[e]) ēt bēllō pāssūs, dūm cōndĕrĕt ūrbĕm,
-- -` - -- -- - -uu -u
īnfērrētquĕ dĕōs Lătĭō, gĕnūs (gĕnŭs) ūndĕ Lătīnŭm,
---u u- uu- -u -u u-u
Ālbānīquĕ pātrēs (pātrēs), ātqu[e] āltāe mōenĭă Rōmāe.
---u -- -` -- -uu --

# What makes the difference ?

- Dictionaries :
a single click opens the dictionary

- Allows to compare two dictionaries

- Integrated server

  → can answer a question from another program

- Expandable (one can add lemmas and/or dictionaries)

- Scansion :
knows a priori if a vowel is short or long

  → not limited to metrical verses

# Principle of the analysis

- Not a list of inflected forms
  → closer to the human mechanism

- Split the word in 2 parts (all possibilities)
- Look for the 2 parts in the proper list
- Check for the compatibility

- 24 100 lemmas in the "main lexicon"
- 58 000 lemmas in the "extended lexicon"

# The figures

- 24 100 lemmas (classical Latin)

- 58 000 lemmas in the extended lexicon (unusual words)

- 140 "paradigms"

- Lewis & Short ≈ 60 000 entries
- Gaffiot 2016 ≈ 60 000 entries
- K.E. Georges ≈ 50 000 entries
- Gérard Jeanneau ≈ 50 000 entries
   Generate automatically the lemmas in the format of Collatinus

7 500 entries still waiting

- Assimilation of the prefix (without an extra entry)

  *obfero ↔ offero*

- Contraction in the declension
  $\breve{i}\breve{i} \to \bar{i}$

  *amasse ↔ amavisse*

# Ordering the analysis

- "Ordering" the analysis
  → a step to desambiguization
  ≠ from a usual POS tagger (no list of forms)


- Statistics on the texts lemmatized by the LASLA (almost 2 000 000 forms)

- Counting the tags, the sequences of 3 tags and each lemma (mapped in the lexicon)

- 91 tags : more than the POS

# TextiColor

- Students are supposed to know a list of words

- The text is prepared with different colors :
  - known
  - not known
  - medieval form

[Scripsit magister Johannes, Veenendaal, Batavia, Anno MMXVII.]

## Navigatio sancti Brendani abbatis

1.1 Sanctus Brendanus, filius Finlocha, nepotis Alti, de genere Eogeni Stagni Len, regionis Mumenensium ortus fuit.
5 Erat vir magne abstinencie et in virtutibus clarus, trium milium fere monachorum pater.

abstinentia, ae, f. = modestia (<modus), continentia, sobrietas
milium: gen. pl. < mille

1.2 Cum esset in suo certamine, in loco qui dicitur Saltus Virtutis Brendani, contigit ut quidam patrum ad eum,
10 quadam vespera, venisset, nomine Barinthus, nepos Neil.

quidam patrum = pater aliquis

1.3 Cumque interrogatus esset multis sermonibus a predicto sancto patre, cepit lacrimari et prostrare se in terram et diucius permanere in oracione. At Sanctus Brendanus erexit
15 illum de terra et osculatus est eum, dicens: "Pater, cur tristiciam habemus in adventu tuo? Nonne ad consolationem nostram venisti? Magis leticiam tu debes fratribus preparare. Indica nobis verbum Dei, atque refice animas nostras de diversis miraculis, quae vidisti in Oceano."

lacrimari: lacrimare | prostrare = prosternere, prostravi, prostratum

20

1.4 Tunc Sanctus Barinthus, expletis his sermonibus sancti Brendani, cepit narrare de quadam insula, dicens: "Filiolus meus Mernoc, procurator pauperum Christi, confugit a facie mea et voluit se esse solitarium. Invenitque insulam iuxta
25 Montem Lapidis, nomine Deliciosam. Post multum vero temporis, nunciatum est mihi quod plures monachos habuisset, et Deus multa mirabilia per illum ostendisset. Itaque perrexi, ut visitassem filiolum meum. Cumque appropinquassem, per trium dierum iter, in occursum mihi
30 festinavit cum fratribus suis. Revelavit enim Dominus sibi adventum meum. Navigantibus nobis in predictam insulam,

ex-plere = h.l. finire

solitarius, i, m. = vir qui solus vivit
delici-osus = quod valde delectat
plures: multos

visita(vi)ssem: visitarem
per: post
sibi: ei/illi

# From quantities to rhythm

- Collatinus knows the quantity of the vowels

  →  easy to determine the tonical accent.

- Can split the syllabes →  hyphenation
  abs·cí·di (abs-cido) ≠ áb·sci·di (ab-scindo)


- Opens the way to the study of rhythmic prose

- Statistics on paroxytons and proparoxytons.

- Rhythm everywhere, not only in *clausulae*
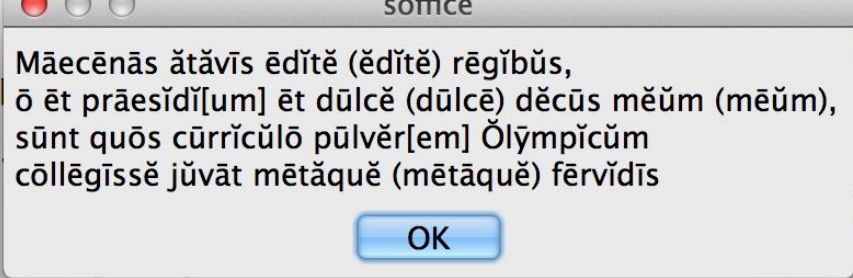
# New feature : probabilistic tagger

- 2nd order hidden Markov model
- Training corpus : LASLA.
- For each sentence, chooses the best sequence of tags and thus the best analysis

- No one knows if a probabilistic tagger works with Latin ("free" order of the words)
- No evaluation of the error rate (not yet)

# Internal server

- Collatinus can answer questions asked from another program. For instance, LibreOffice.

# Internal server (2)

- A "standard" tagger with the data from LASLA
  - List of forms, with lemma and analysis
  - Exact number of occurrences
  - Tagset and statistics on sequences of tags

- Problem of the forms that have not been seen before : Ask Collatinus !

  → Answers a LASLA-like analysis

# LASLA-Tagger

- Supervised lemmatization.
- Edition/correction of the result of the tagger.

  → Hope to reach zero error !

- Exact number of occurrences. For "patres" :
  253 voc. 110 nom. and 75 acc.
  Collatinus expects 6 voc. 200 nom. & 320 acc

# LASLA-Tagger : screen-shot

# What is going on ?

***Praelector:***

- "Grammatical assistant" to help pupils in reading Latin / understanding the sentence.

  → "Translation" in French of the sentence

- Medieval spelling :
  e for <span style="color:red">ae</span>, o for <span style="color:red">au</span>, f for <span style="color:red">ph</span>, etc...
  Reduce the medieval forms and the classical words to the same "phonetic" form.

# Future plans

- Coupling of the probabilistic tagger with the gramatical analysis of the sentence

- Statistics on rhythmic prose (cursus)

- Looking for verses in prose

- "Universal" output format to gather all the information about the text (lemmas, analysis, scanned and accented forms...)

# Suggestions

- Yves.Ouvrard@collatinus.org

- Philippe.Verkerk@univ-lille1.fr

Thank you for your attention

- http://outils.biblissima.fr/en/collatinus/

- http://outils.biblissima.fr/en/collatinus-web/

- For ancient Greek : http://outils.biblissima.fr/en/eulexis/