

Open Philology Workshop

8.-10. August 2013, Leipzig University

An Introduction to a Latin Language Toolkit

Prometheus*

* working title

Overview

- Background
- Principles and Goals
- Understanding Latin
 - Tokenization
 - Morphological Analysis
 - Syntactical Analysis - the system's internals
 - Evaluating Caesar
 - Upcoming Improvements and Challenges
- Linking with Data-driven parsers - the system as a Post-Processing Unit
- Creating Latin
 - A framework for Latin studies
 - Dynamic e-learning tools
- The Future - cooperating with the research community

Background

□ Who are we?

- no experience in computer linguistics
- Study of Latin Philology, Classical History, History (teacher-certificate)
- Experience in teaching Latin

□ Ideas and Motivation

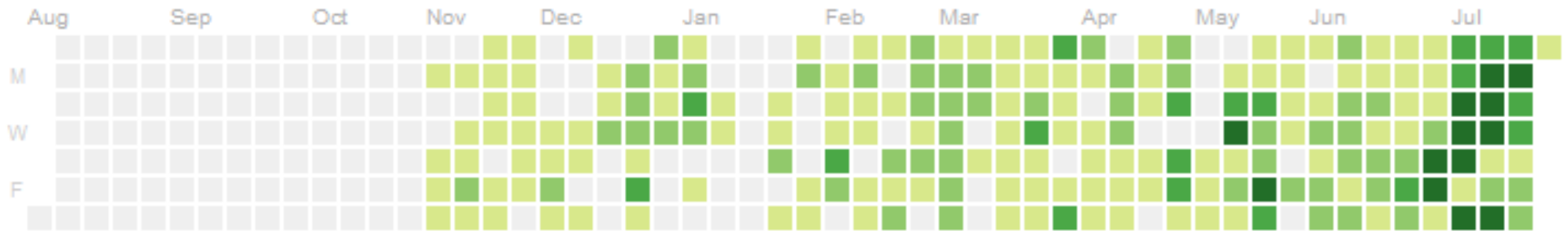
- Author specific exercises
- Immediate feedback

□ Working conditions

- Private project
- Built from scratch

Timeline

Contributions



? Summary of Pull Requests, issues opened and commits. [Learn more.](#)

Less More

1,273 Total

Aug 04 2012 - Aug 04 2013

Year of Contributions

55 days

June 11 - August 04

Longest Streak

55 days

June 11 - August 04

Current Streak

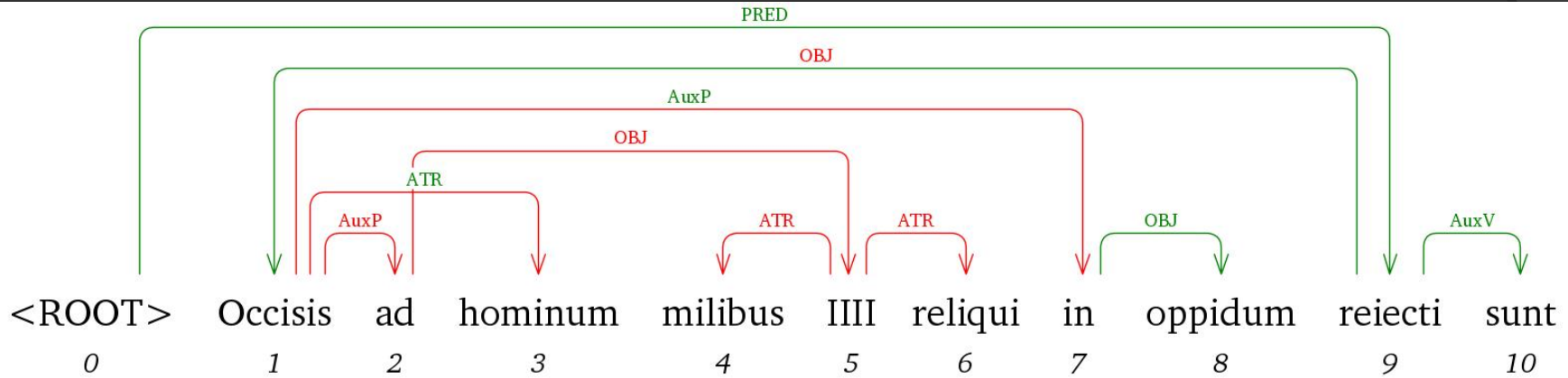
Principles

- ▮ Texts are written to be understood
- ▮ No special training required
- ▮ Semantics are not deciding
- ▮ Non-projectivity is nothing special

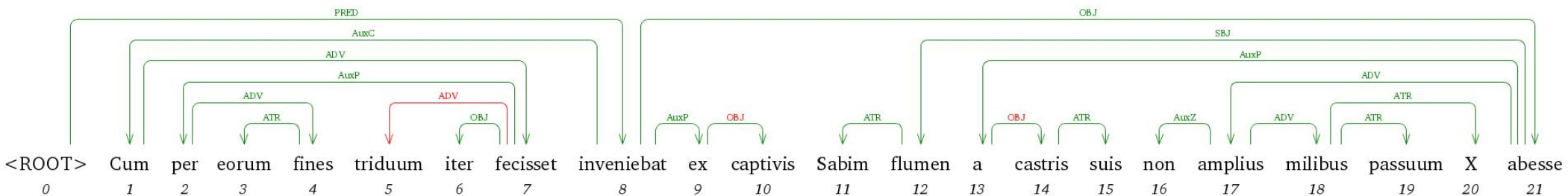
Requirements

- ▮ Two-way-engine human style
- ▮ Modular, extendable and flexible
- ▮ Easy recalibration
- ▮ No results are better than malformed results

malformed result:



acceptable result:



Requirements

- ▮ Two-way-engine human style
- ▮ Modular, extendable and flexible
- ▮ Easy recalibration
- ▮ No results are better than malformed results
- ▮ One language done well
- ▮ “real-time” computation speed

=> object oriented programming

Understanding Latin

-

the parsing sequence

Tokenization

Example sentence (Caes. BG 2,9,5):

Ubi neutri transeundi initium faciunt, secundiore equitum proelio nostris Caesar suos in castra reduxit.

„As neither side started to cross (the river), Caesar led his troops back to the camp, after a cavalry battle, which favoured our men.”

Tokenization

Example sentence (Caes. BG 2,9,5):

Ubi neutri transeundi initium faciunt ,
secundiore equitum proelio nostris Caesar
suos in castra reduxit .

Morphology

- Analysis of words
 - Morphem-specification (stem, thematic-vowels, tempus/modus signs/marker, ending)
- Database look-up
 - Stem dictionary

Database entries

- ▮ „castra“ will return 4 forms:
 - nom. pl. n. – O. Decl
 - acc. pl. n. – O. Decl
 - voc. pl. n. – O. Decl
 - Imperative, pr. 2. p. sg. – A. Conj.

Morphology

- Analysis of words
 - Morphem-specification (stem, thematic-vowels, tempus/modus signs/marker, ending)
- Database look-up
- Phonology and metrical info
- Didactic aspects (patterns!)
 - ama – ba – t

Syntax - process overview

- ▮ Clause finding process
- ▮ Clause analysis
- ▮ Clause rating
- ▮ Clause merging

Clause Finding

Ubi neutri transeundi initium **faciunt** ,
secundiore equitum **proelio** nostris Caesar
suos in **castra reduxit**.

Safe Subjunction / Verb / Comma

Unsafe Verb

Clause Finding

Clause 1 (subordinate):

(**Ubi** neutri transeundi initium **faciunt** ,)

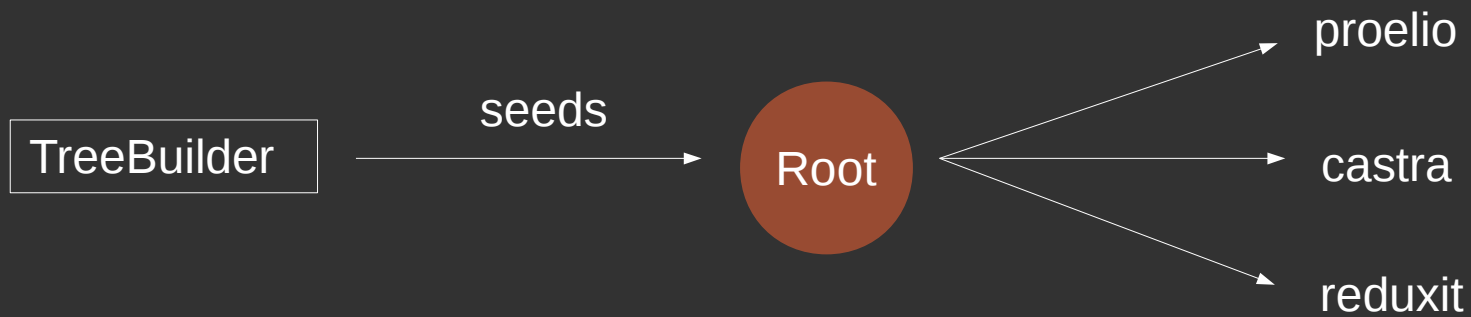
Clause 2 (main):

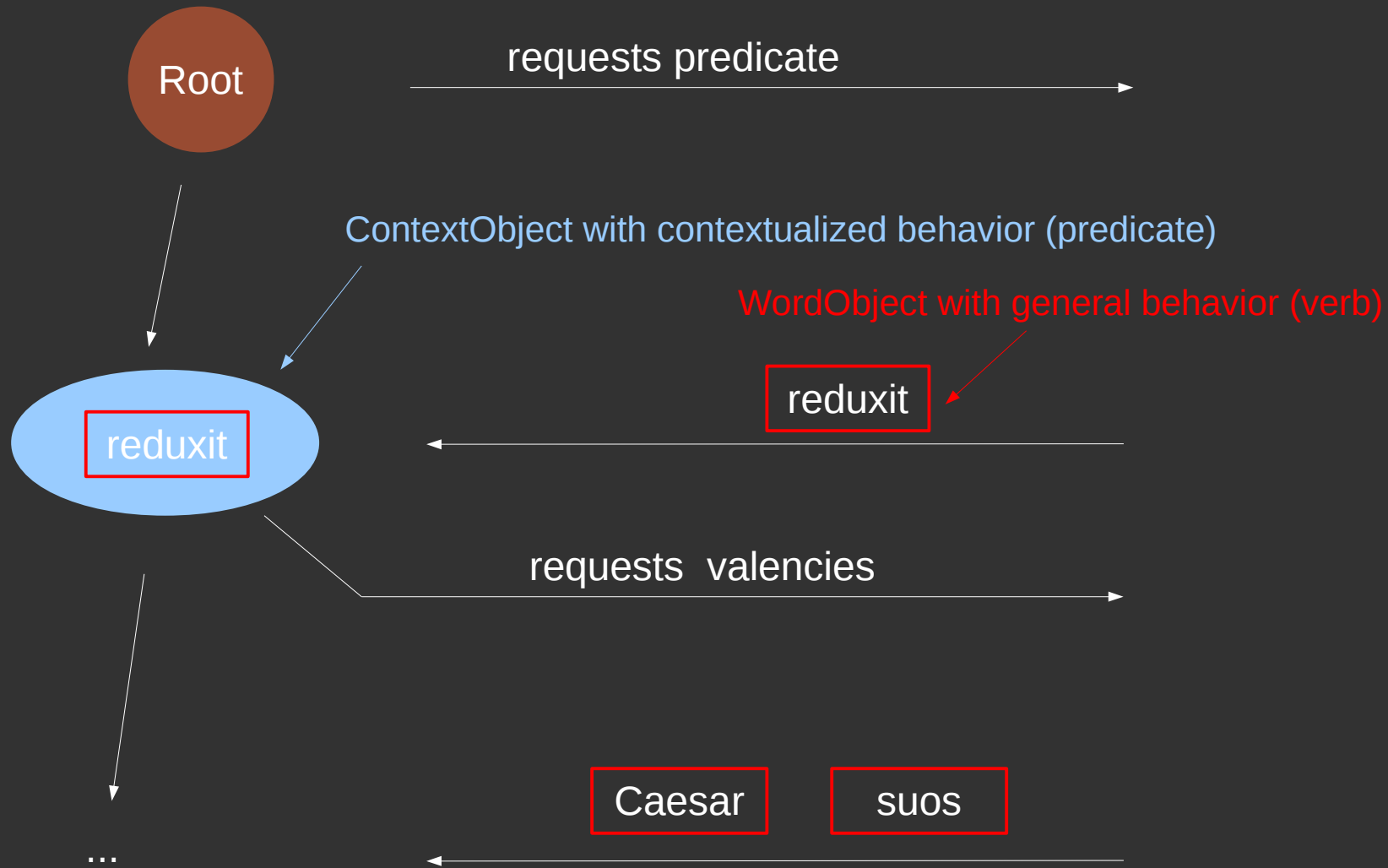
(secundiore equitum **proelio** nostris Caesar suos in **castra reduxit.**)

Safe Subjunction / Verb / Comma

Unsafe Verb

Clause analysis (tree-building)





UNUSED

Constraint Behavior

amicus **militis** salutat.

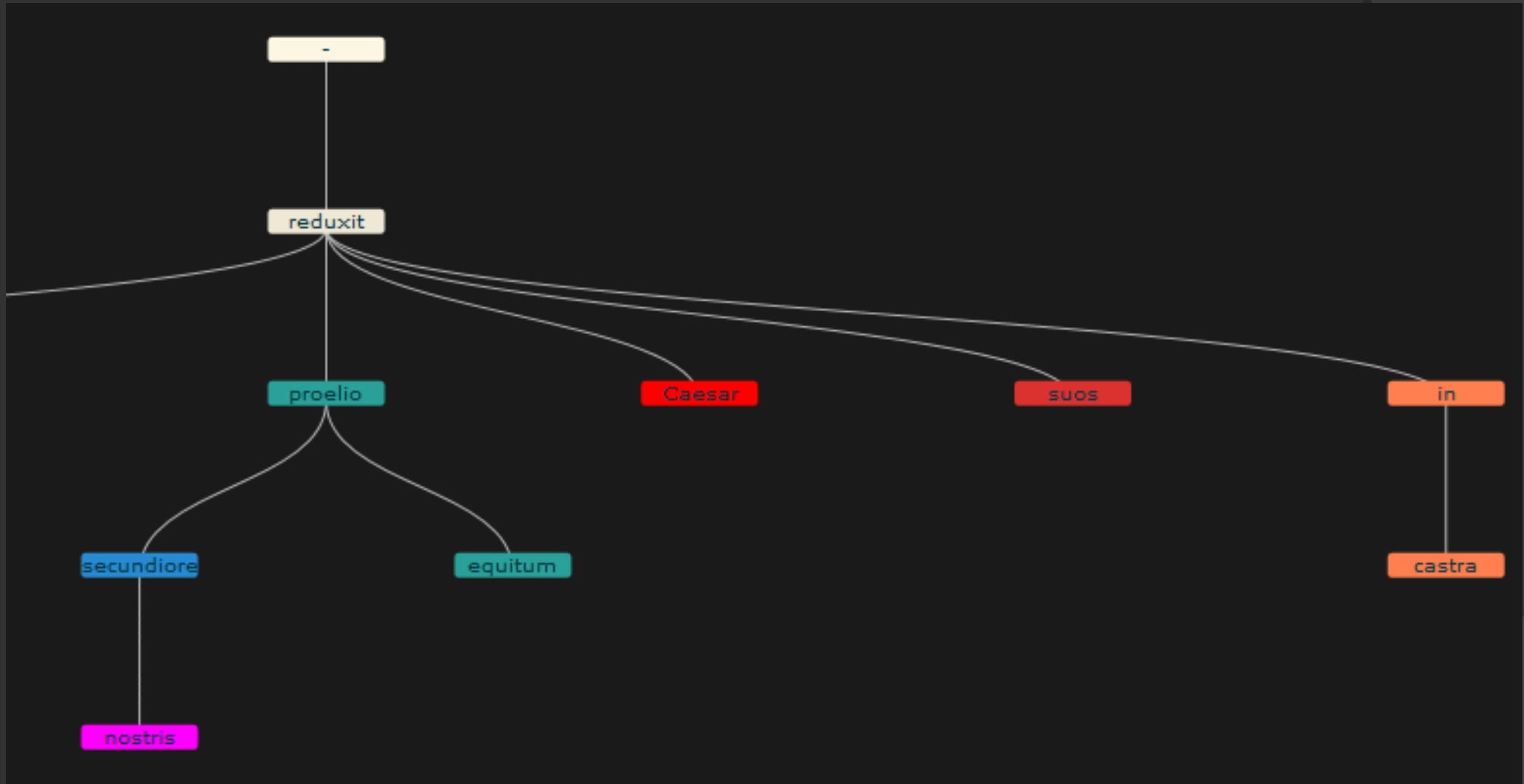
(„The soldier’s friend“)

amicus **miles** **Caesaris** pugnat.

(„Caesar’s friendly soldier“)

SBJ / **ATR (congruence)** / **ATR (genetive)**

Full clause leaf



Ratings

Et id quidem nemini video Graecorum adhuc contigisse

134

Awards

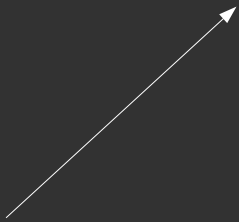
50 for 1x nominalized neutrum pl
30 for 1x lone infinitive is regens
25 for 1x subj acc
25 for 1x close to regens
20 for 1x root encapsulation
20 for 1x object in valency group
20 for 1x is nested
19 for 1x root dependent
10 for 1x indirect object
5 for 1x animated

Penalties

-20 for 1x has neutrum subject
-70 for 1x adjective is nominalized

Clause Merging

secundiore equitum proelio nostris Caesar suos in
castra **reduxit**



Ubi neutri transeundi initium faciunt

Problems

Eo cum de improviseo celeriusque omnium opinione venisset, Remi, qui proximi Galliae ex Belgis sunt, ad eum legatos Iccium et Andebrogium, primos civitatis, miserunt, qui dicerent se suaque omnia in fidem **atque** potestatem populi Romani **permittere**, ne**que** se cum reliquis Belgis **consensisse** ne**que** contra populum Romanum **coniurasse**, **paratosque esse et** obsides **dare et** imperata **facere et** oppidis **recipere et** frumento ceterisque rebus **iuvare**;

Permittere => consensisse => coniurasse => paratos esse => dare => facere => recipere => iuvare

Evaluation

Caesar LDT annotation*:

LAS

UAS

LA

84,7 %

91,0 %

89,7 %

**preliminary results (85% of all tokens, without punctuation)*

40% of the parsed sentences at 100 %

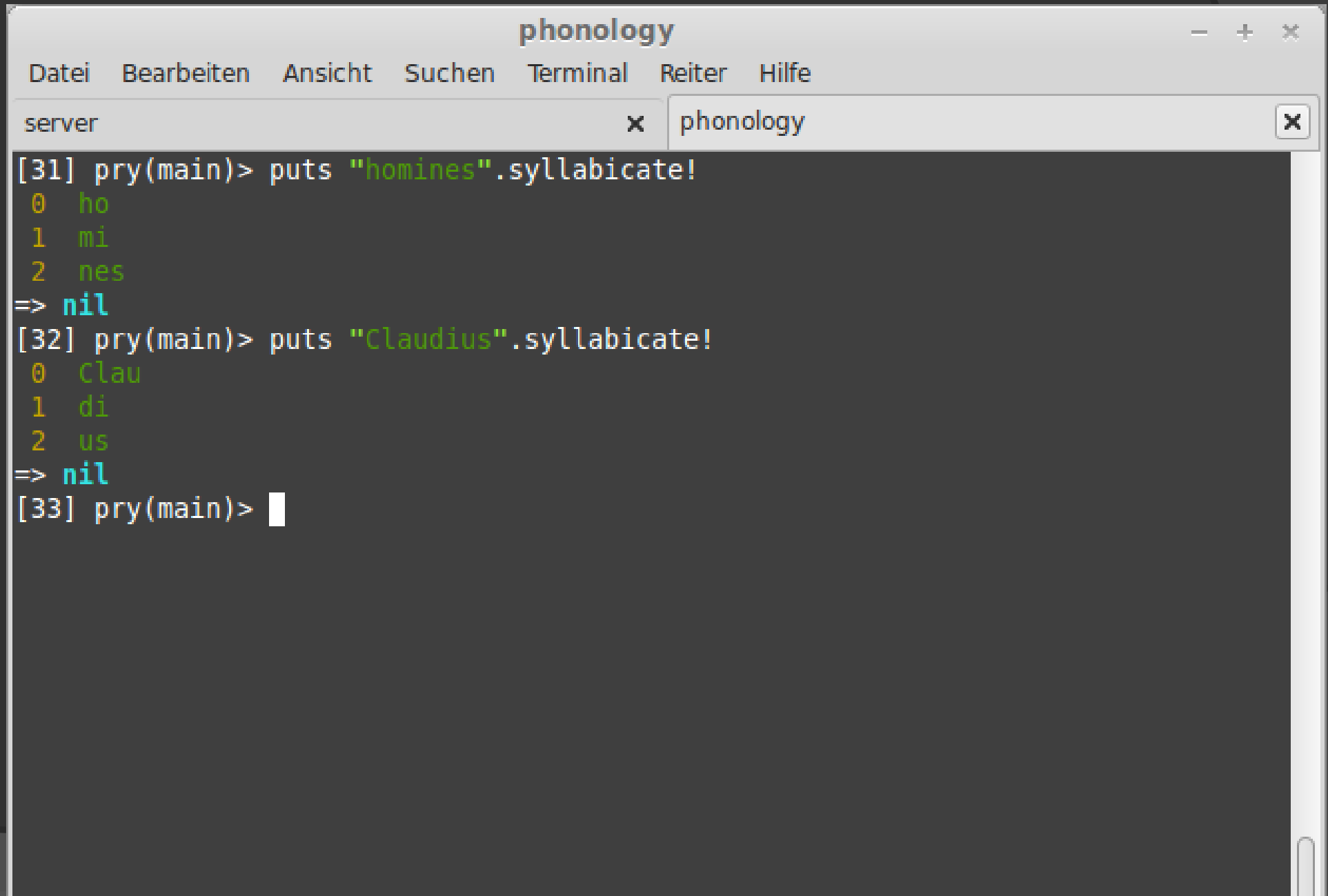
Challenges & Improvements

- ▮ Poetry

Basic-Tools: Phonology

```
phonology
Datei Bearbeiten Ansicht Suchen Terminal Reiter Hilfe
server x phonology x
[26] pry(main)> puts "homines".phonologize!
0 h consonant glottal spirant
1 o vowel
2 m consonant labial nasal sonant
3 i vowel
4 n consonant dental nasal sonant velar
5 e vowel
6 s consonant dental spirant
=> nil
[27] pry(main)> puts "Claudius".phonologize!
0 C consonant muta tenuis velar
1 l consonant dental liquida sonant
2 au diphtong
3 d consonant dental media muta
4 i vowel
5 u vowel
6 s consonant dental spirant
=> nil
[28] pry(main)> █
```

Basic-Tool: Syllables

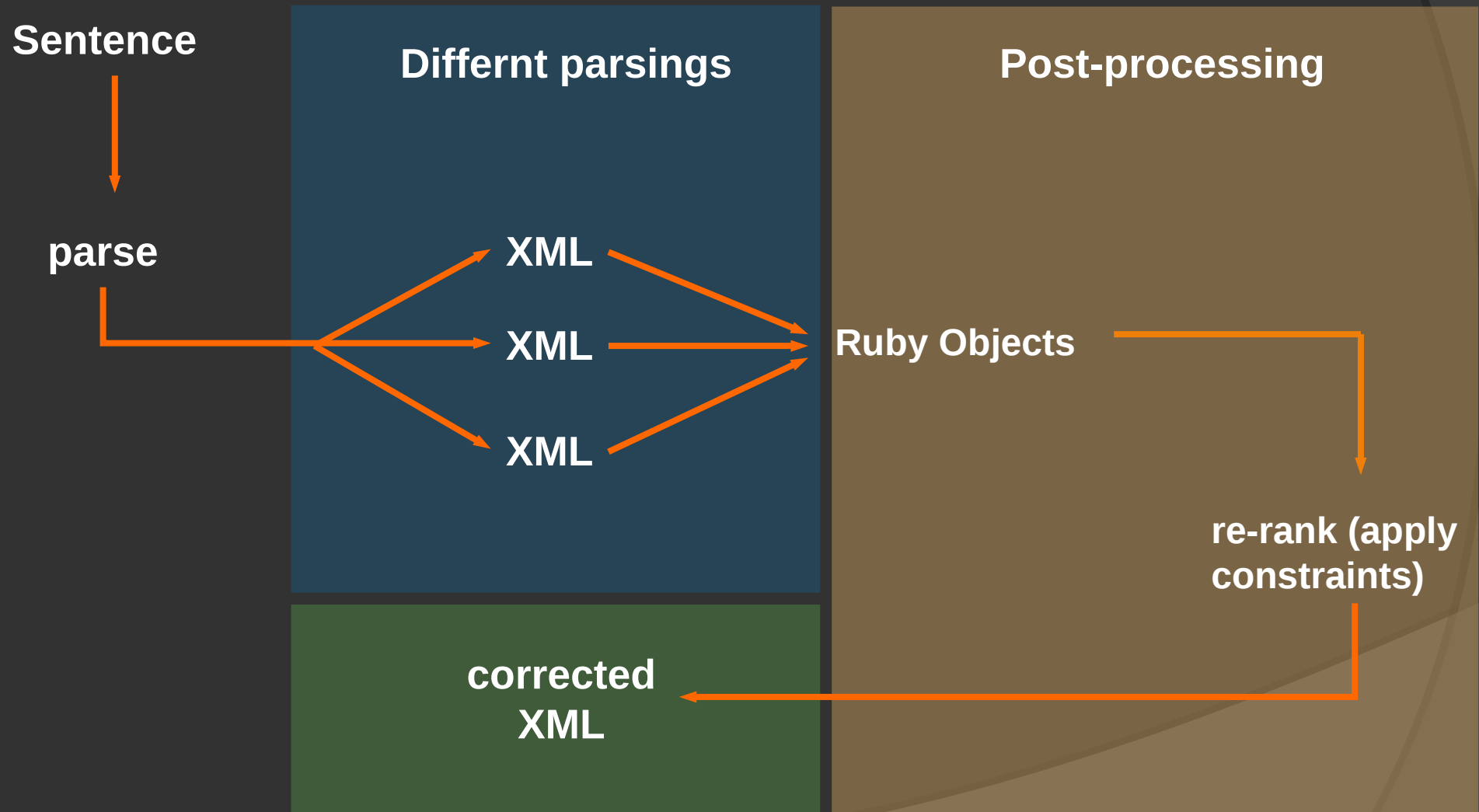


```
phonology
Datei Bearbeiten Ansicht Suchen Terminal Reiter Hilfe
server x phonology x
[31] pry(main)> puts "homines".syllabicate!
0 ho
1 mi
2 nes
=> nil
[32] pry(main)> puts "Claudius".syllabicate!
0 Clau
1 di
2 us
=> nil
[33] pry(main)> █
```

Challenges & Improvements

- ▮ Poetry
- ▮ Semantics
- ▮ Dictionary/Database

Post-Processing Data-driven parsings



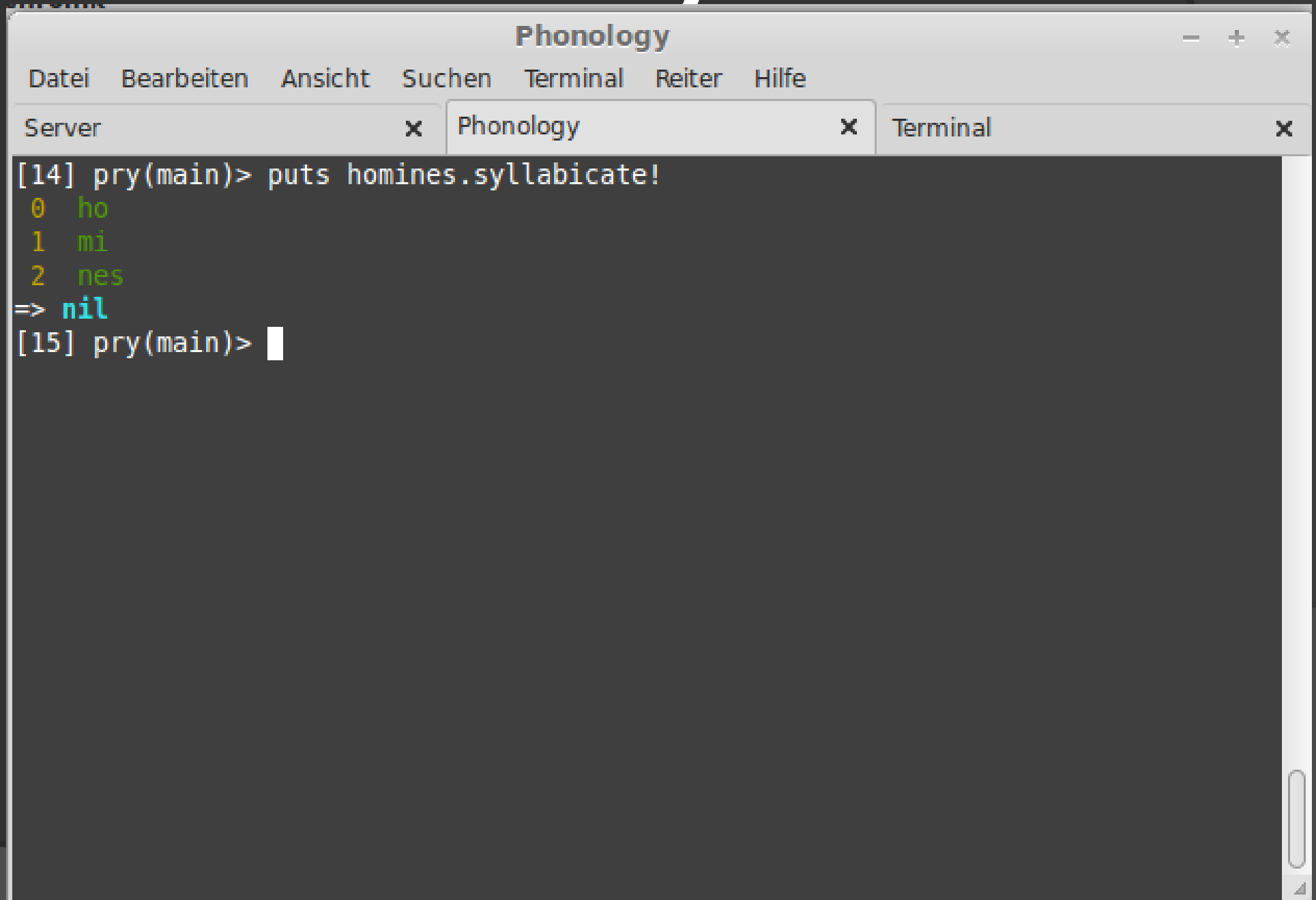
**Thanks for your
patience and
attention!**

to discuss...

Basic-Tools: Phonology

```
phonology
Datei Bearbeiten Ansicht Suchen Terminal Reiter Hilfe
server x phonology x
[40] pry(main)> "homines".tokenize!(check: true)
=> [Word token: homines
      Forms: 1: 1m2 CDe | 2: 5m2 CDe | 3: 4m2 CDe
      checked]
[41] pry(main)> "homines".tokenize!(check: true).first
=> Word token: homines
      Forms: 1: 1m2 CDe | 2: 5m2 CDe | 3: 4m2 CDe
      checked
[42] pry(main)> "homines".tokenize!(check: true).first.use(1)
=> homines 1m2 CDe
[43] pry(main)> homines = _
=> homines 1m2 CDe
[44] pry(main)> puts homines.phonologize!
0 h consonant glottal spirant
1 o vowel
2 m consonant labial nasal sonant
3 i vowel
4 n consonant dental nasal sonant velar
5 e vowel
6 s consonant dental spirant
=> nil
[45] pry(main)> █
```

Basic-Tool: Syllables



The image shows a window titled "Phonology" with a menu bar containing "Datei", "Bearbeiten", "Ansicht", "Suchen", "Terminal", "Reiter", and "Hilfe". Below the menu bar are three tabs: "Server", "Phonology", and "Terminal". The "Phonology" tab is active and contains a terminal window with the following text:

```
[14] pry(main)> puts homines.syllabicate!  
0 ho  
1 mi  
2 nes  
=> nil  
[15] pry(main)> █
```

Clause Finding

Caes. BG 2,34

...quem cum legione una miserat ad
Venetos, Venellos, Osismos, Coriosolitas,
Esvivos, Aulercos, Redones...

Parsing Poetry

Ovid, Meta., 10,529

Capta viri forma non iam Cythereia curat | litora

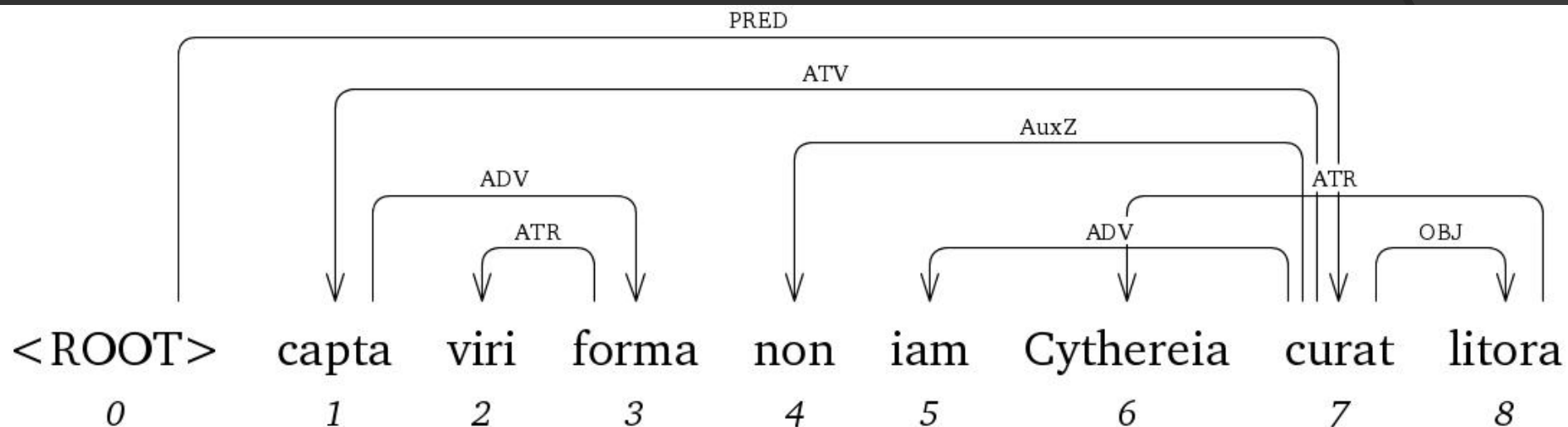
1: 1f1 PPP
2: 1n2 PPP
3: 4n2 PPP
4: 5n2 PPP
5: 5f1 PPP
6: 6f1 PPP
7: 1n2 PPP
8: 4n2 PPP
9: 5n2 PPP

1: 1f1 ADe
2: 5f1 ADe
3: 6f1 ADe
4: 1f1 AOD
5: 5f1 AOD
6: 6f1 AOD
7: 1n2 AOD
8: 4n2 AOD
9: 5n2 AOD
10: 2,1.i.a.pr ACo

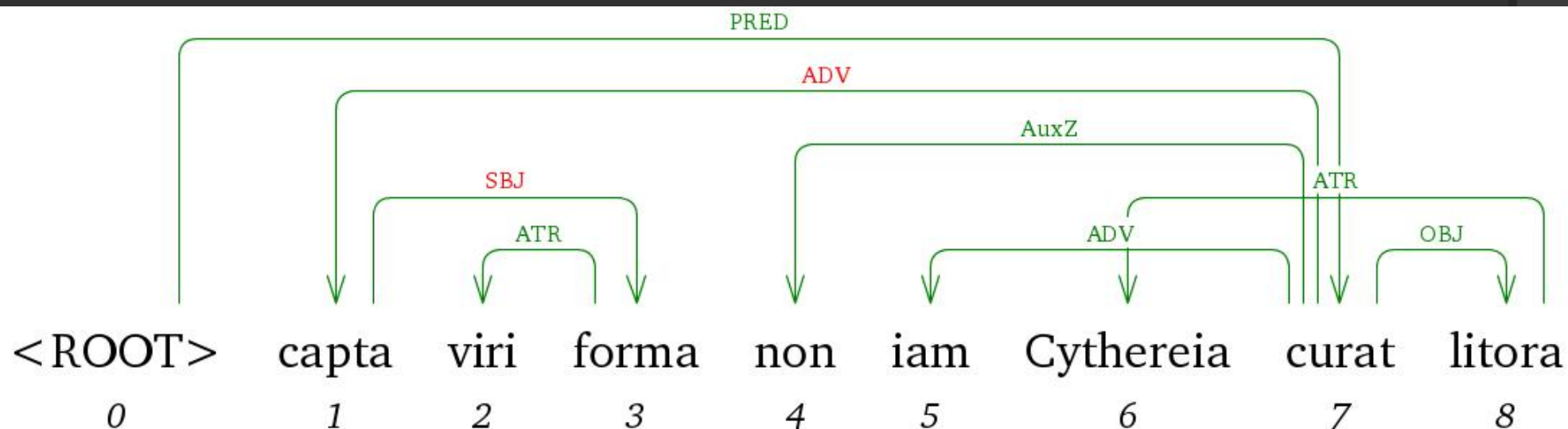
1: 1f1 AOD
2: 5f1 AOD
3: 6f1 AOD
4: 1n2 AOD
5: 4n2 AOD
6: 5n2 AOD

1: 1n2 CDe
2: 5n2 CDe
3: 4n2 CDe

Gold-Standard:



Result:



Parsing Poetry

Căptă vîrî fōrmā nōn iām Cūthēreîă cūrăt | lītōră,

1: 1f1 PPP

2: 1n2 PPP

3: 4n2 PPP

4: 5n2 PPP

5: 5f1 PPP

~~6: 6f1 PPP~~

7: 1n2 PPP

8: 4n2 PPP

9: 5n2 PPP

~~1: 1f1 ADe~~

~~2: 5f1 ADe~~

3: 6f1 ADe

~~4: 1f1 AOD~~

~~5: 5f1 AOD~~

6: 6f1 AOD

~~7: 1n2 AOD~~

~~8: 4n2 AOD~~

~~9: 5n2 AOD~~

10: 2,1.i.a.pr ACo

1: 1f1 AOD

2: 5f1 AOD

~~3: 6f1 AOD~~

4: 1n2 AOD

5: 4n2 AOD

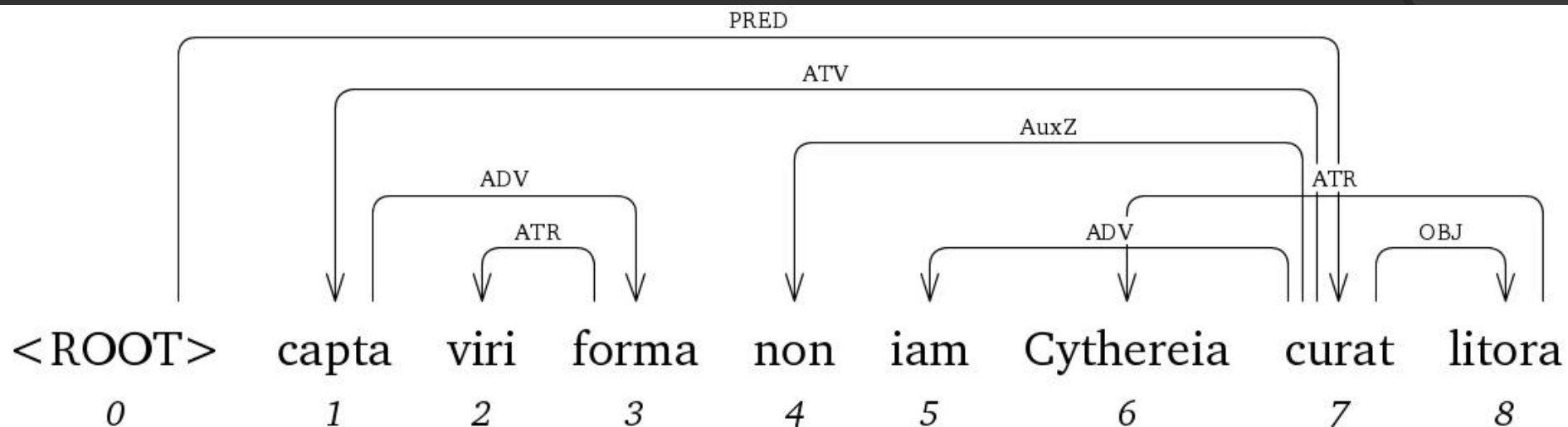
6: 5n2 AOD

1: 1n2 CDe

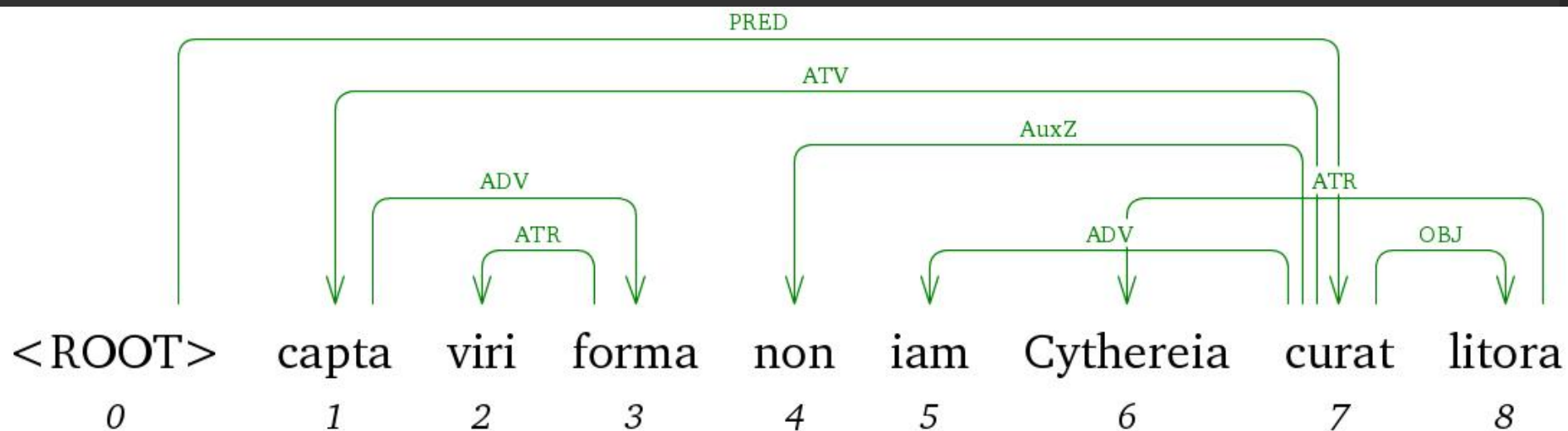
2: 5n2 CDe

3: 4n2 CDe

Gold-Standard:



New result:



Performance/computational speed

Cum esset Caesar in citeriore Gallia in hibernis ita, uti supra demonstravimus, crebri ad eum rumores adferebantur litterisque item Labieni certior fiebat omnes Belgas, quam tertiam esse Galliae partem dixeramus, contra populum Romanum coniurare obsidesque inter se dare.

45 tokens => 7.3 seconds

Dat negotium Senonibus reliquisque Gallis, qui finitimi Belgis erant, uti ea, quae apud eos gerantur, cognoscant seque de his rebus certiozem faciant. **29 tokens => 1.7 seconds**

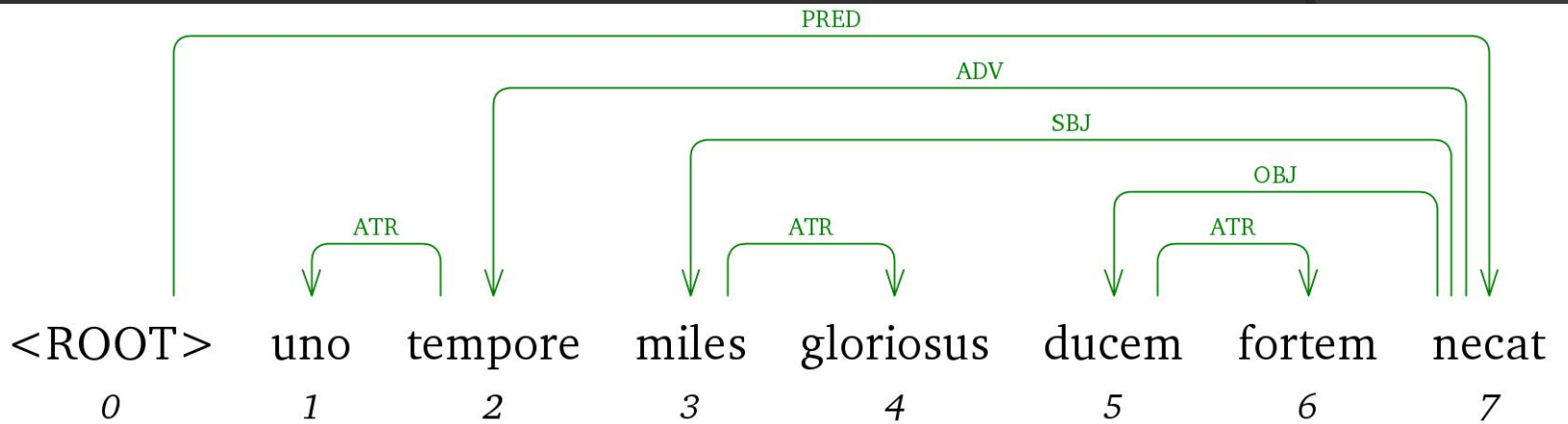
His nuntiis litterisque commotus Caesar duas legiones in citeriore Gallia novas conscripsit et in ita aestate in ulteriorem Galliam, qui deduceret, Q. Pedium legatum misit.

28 tokens => 4.4 seconds

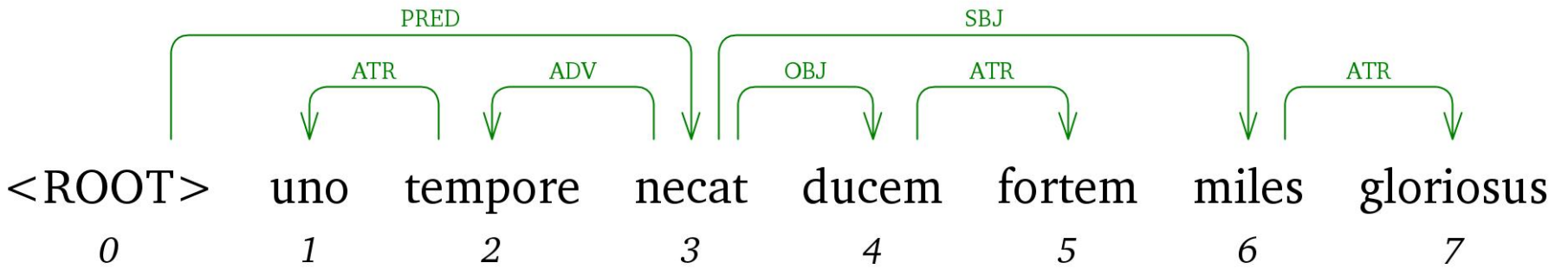
Performance/computational speed

- Intel i5-3570k (3.4GHz), 8GB RAM (4 cores)
- All Caesar-tokens (~ 1400):
2 min. 20 sec
- ~ 500 tokens per minute

SOV:



OSV:



Non-projectivity:

SOV:

