SunoikisisDC Summer 2024

Session 9

# Treebanking and annotation to study reception of ancient language and text

Marja Vierros (U Helsinki)

Alek Keersmaekers (KU Leuven)

Gabriel Bodard (U London)

# Julia Balbilla poem 1

Ἰουλίας Βαλ<β>ίλλης·
ὅτε ἤκουσε τοῦ Μέμνο<νο>ς
ὁ Σεβαστὸς Ἀδριανός.

Μέμνονα πυνθανόμαν Αἰγύπτιον, ἀλίω αὔγαι
    αἰθόμενον, φώνην Θηβαΐ<κ>ω ʼπυ λίθω.
Ἀδρίανον δʼ ἐσίδων, τὸν παμβασίληα, πρὶν αὔγας
    ἀελίω χαίρην εἶπέ <ϝ>οι ὡς δύνατον.
Τίταν δʼ ὅττʼ ἐλάων λεύκοισι διʼ αἴθερος ἵπποις
    ἐγ̣ὶ σκίαι ὠράων δεύτερον ἦχε μέτρον,
ὡς χάλκοιο τύπεντ[ο]ς ἵη Μέμνων πάλιν αὔδαν
    ὀξύτονον· χαίρω[ν κ]αὶ τρίτον ἄχον ἵη.
κοίρανος Ἀδρίανο[ς τότʼ ἄ]λις δʼ ἀσπάσσατο καῦτος
    Μέμνονα, κἀν [στά]λαι κάλ̣λι[π]εν ὀψ[ι]γόνοις
γρόππατα σαμαίνο[ν]τά τʼ ὅσʼ εὔϊδε κὤσσʼ ἐσάκουσε,
    δῆλον παῖσι δʼ ἔγε[ν]τʼ ὥς <ϝ>ε φίλισι θέοι.

# Julia Balbilla poem 2

ὅτε σὺν τῇ Σεβαστῇ Σαβείνη-
   ι ἐγενόμην παρὰ τῷ Μέμνονι.
Αὔως καὶ γεράρω, Μέμνον, πάϊ Τιθώνοιο,
   Θηβάας θάσσων ἄντα Δίος πόλιος,
ἢ Ἀμένωθ, βασίλευ Αἰγύπτιε, τὼς ἐνέποισιν
   ἵρηες μύθων τῶν παλάων ἴδριες,
χαῖρε, καὶ αὐδάσαις πρόφρων ἀσπάσδε[ο κ]αὔτ[αν]
   τὰν σέμναν ἄλοχον κοιράνω Ἀδριάνω.

γλῶσσαν μέν τοι τμᾶξε [κ]αὶ ὤατα βάρβαρος ἄνηρ,
   Καμβύσαις ἄθεος· τῶ ῥα λύγρω θαγάτω
δῶκέν τοι ποίναν τὤτωι ἄκ[ρω] ἄορι πλάγεις
   τῷ νήλας Ἆπιν κάκτανε τὸν θέϊον.
ἀλλ᾽ ἔγω οὐ δοκίμωμι σέθεν τόδ᾽ ὄλεσθ᾽ ἂν ἄγαλμα,
   ψύχαν δ᾽ ἀθανάταν λοίπον ἔσωσα νόω.

εὐσέβεες γὰρ ἔμοι γένεται πάπποι τ᾽ ἐγένοντο,
   Βάλβιλλός τ᾽ ὁ σόφος κ᾽ Ἀντίοχος βασίλευς,
Βάλβιλλος γενέταις μᾶτρος βασιλήϊδος ἄμμας,
   τῶ πάτερος δὲ πάτηρ Ἀντίοχος βασίλευς·
κήνων ἐκ γενέας κἄγω λόχον αἶμα τὸ κᾶλον,
   Βαλβίλλας δ᾽ ἔμεθεν γρόπτα τάδ᾽ εὐσέβεος.

# Julia Balbilla poem 3

ὄτε τῇ πρώτῃ ἡμέρᾳ οὐκ ἀ-
κούσαμεν τοῦ Μέμνονος.


χθίσδον μέν Μέμνων σίγαις ἀπε[δέξατ᾽ ἀκ]οίτα[ν](?),
  ὡς πάλιν ἀ κάλα τυῖδε Σάβιννα μό[λοι].
τέρπει γάρ σ᾽ ἐράτα μόρφα βασιλήϊδος ἄμμας,
  ἐλθοίσαι δ᾽ [α]ὔται θήϊον ἄχον ἴη,
μὴ καί τοι βασίλευς κοτέσῃ, τό νυ δᾶρον ἀτά[ρβης]
  τὰν σέμναν κατέχες κουριδίαν ἄλοχον.
κὠ Μέμνων τρέσσαις μεγάλω μένος Ἀδρι[άνοιο]
  ἐξαπίνας αὔδασ᾽, ἀ δ᾽ ὀΐοισ᾽ ἐχάρη.

# Julia Balbilla poem 4

ἔκλυον αὐδήσαντος ἔγω ’πυ λίθω Βάλβιλλα
    φώνα<ς> τᾶς θείας Μέμνονος ἢ Φαμένωθ.
ἦλθον ὔμοι δ’ ἐράται βασιλήϊδι τυῖδε Σαβίννᾳ,
    ὤρας δὲ πρώτας ἄλιος ἦχε δρόμος.
κοιράνω{ι} Ἀδριάνω πέμπτῳ δεκότῳ δ’ ἐ/νιαύτῳ,
    <φῶτ>α δ’ ἔχεσκε<ν> Ἄθυρ εἴκοσι / καὶ πέσυρα.

- Encoding texts in EpiDoc
- Dialect analysis, text emendation
- Epigraphic + archaeological metadata
- Encoding/analysing metre
- Prosopography/genealogy, religion
- Image annotation
- Translation alignment
- Morphosyntactic annotation / syntax analysis
- Reception of archaic poetry in imperial literature

# Treebanking & comparing syntax Case study: Julia  Balbilla

Marja Vierros, University of Helsinki

# Outline today

- Very short intro to treebanking *aka* linguistic annotation
- Very short intro to PapyGreek platform for annotating papyri (and inscriptions)
- Demonstration of treebanking Balbilla's poem 1
  - For a full tutorial of treebanking, see the Sunoikisis Session 9 page (exercise), or
  - Other tutorial listed in https://wiki.digitalclassicist.org/Treebanking
  - Remember to read the guidelines, too!
- Comparing syntactic structures of Balbilla to Sappho (only preliminary, non-machine-facilitated glimpse)

# What is linguistic annotation?

- many different types of linguistic information can be added to a (digital) text
  - WORD LEVEL
    - lemma annotation (lemma = basic form of a word, the 'dictionary entry')
    - morphological annotation (inflectional morphology like case, number, gender etc.)
    - POS annotation (word classes i.e. parts-of-speech)
  - SYNTACTIC ANNOTATION
    - different grammar formalisms: most important here: *dependency* and *constituent* structures
  - SEMANTIC ANNOTATION
    - Semantic classes based on meanings of words (or phrases or sentences), e.g. concrete / abstract nouns; animate / inanimate etc.
    - Annotating concepts (e.g. people, places, organizations, products or topics)

# What is a treebank?

- tree structure (tree diagram) is a graphical way of representing a **hierarchical** structure, e.g. the syntax of a sentence

- Treebank = parsed text corpus that annotates syntactic (or semantic) sentence structure
  - Usually includes lemmatization and morphological annotation
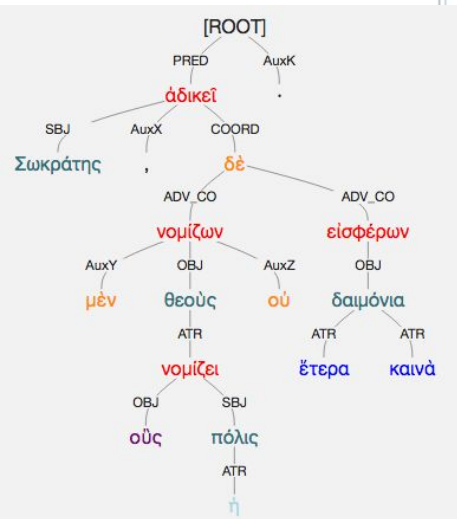


*Constituency tree*

*Dependency tree*

# Treebanks of Ancient Greek and Latin (=Dependency Trees)

- The "Index Thomisticus Treebank" (inspired by the Prague Dependency Treebank)
- The "Perseus family" (AGLDT, Pedalion, Gorman Trees, PapyGreek …)
- PROIEL = Pragmatic Resources in Old Indo-European Languages

- **Universal Dependencies** (UD)

# Underlying XML of a dependency tree (AGDT)



```
<sentence id="3"
  document_id="http://perseids.org/cts5/ne                    :greekLit:tlg003...        ...as-grc2"
  subdoc="1.1.1-1.1.20"
  span="">
  <word id="1" form="ἀδικεῖ" lemma="ἀδικέω" postag="v3spia---" relation="PRED" head="0"/>
  <word id="2" form="Σωκράτης" lemma="Σωκράτης" postag="n-s---mn-" relation="SBJ" head="1"/>
  <word id="3" form="οὓς" lemma="ὅς" postag="p-p---ma-" relation="_BJ" head="7"/>
  <word id="4" form="μὲν" lemma="μέν" postag="d--------"        AuxY" head="10"/>
  <word id="5" form="ἡ" lemma="ὁ" postag="l-s---fn-"                   "6"/>
  <word id="6" form="πόλις" lemma="πόλις" postag="n-               BJ" head="7"/>
  <word id="7" form="νομίζει" lemma="νομίζω" postag=               "ATR" head="8"/>
  <word id="8" form="θεοὺς" lemma="θεός" postag="n-p-             OBJ" head="10"/>
  <word id="9" form="οὐ" lemma="οὐ" postag="d--------"            AuxZ" head="10"/>
  <word id="10" form="νομίζων" lemma="νομίζω" postag="v-sppamn-" relation="ADV_CO" head="13"/>
  <word id="11" form="," lemma="punc1" postag="u--------" relation="AuxX" head="1"/>
  <word id="12" form="ἕτερα" lemma="ἕτερος" postag="a-p---na-" relation="ATR" head="15"/>
  <word id="13" form="δὲ" lemma="δέ" postag="d--------" relation="COORD" head="1"/>
  <word id="14" form="καινὰ" lemma="καινός" postag="a-p---na-" relation="ATR" head="15"/>
  <word id="15" form="δαιμόνια" lemma="δαιμόνιον" postag="n-p---na-" relation="OBJ" head="16"/>
  <word id="16" form="εἰσφέρων" lemma="εἰσφέρω" postag="v-sppamn-" relation="ADV_CO" head="13"/>
  <word id="17" form="·" lemma="punc1" postag="u--------" relation="AuxK" head="0"/>
</sentence>
```

POSTAG: all morphological information

HEAD: which word is governing this one

RELATION: syntactic role

# postag

- nine place string including the morphological information

**ἀδικεῖ**

postag="v3spia---"

1: verb
2: 3<sup>rd</sup> person
3: singular
4: present
5: indicative
6: active
7: -
8: -
9: -

1: part-of-speech
2: person
3: number
4: tense
5: mood
6: voice
7: gender
8: case
9: degree

**πόλις**

postag="n-s---fn-"

1: noun
2: -
3: singular
4: -
5: -
6: -
7: feminine
8: nominative
9: -

Find the lettercodes here: https://github.com/gcelano/LemmatizedAncientGreekXML

https://papygreek.com/

# PapyGreek

A platform for the linguistic study of Greek papyri,
including a **grammar**, **annotated texts**, and a **search tool**.



Grammar not public yet

PapyGreek Treebanks 3.0:
https://zenodo.org/records/8428823

Whole corpus of papyri for phonology and
morphology; treebanked corpus for syntax

# PapyGreek platform for treebanking papyri (and inscriptions)

- Greek Documentary Papyri all there (synchronized with Papyri.info)
- TEI EpiDoc XML preprocessed to get annotatable texts
  - Papyri and inscriptions are fragmentary, contain markup, and cannot be annotated as is
- –> works also for inscriptions (although some differences in markup)
- Arethusa annotation tool implemented within PapyGreek
  - Arethusa originally exists in Perseids' page (Tufts University): https://sosol.perseids.org/sosol/

Demo on Julia Balbilla's inscribed poems to follow

# Julia Balbilla's poem 1 (beginning)

Ἰουλίας Βαλ<β>ίλλης·
    ὅτε ἤκουσε τοῦ Μέμνο<νο>ς
    ὁ Σεβαστὸς Ἀδριανός.

Μέμνονα πυνθανόμαν Αἰγύπτιον, ἀλίω αὔγαι
    αἰθόμενον, φώνην Θηβαΐ<κ>ω 'πυ λίθω.
Ἀδρίανον δ' ἐσίδων, τὸν παμβασίληα, πρὶν αὔγας
    ἀελίω χαίρην εἶπέ <ϝ>οι ὡς δύνατον.

[Composed] by Julia Balbilla,
When the august Hadrian
heard Memnon.

I've heard tell that Memnon the Egyptian, warmed by the sun's rays,
  utters a loud sound from the Theban stone.
Seeing Hadrian, the greatest of kings, and before greeting
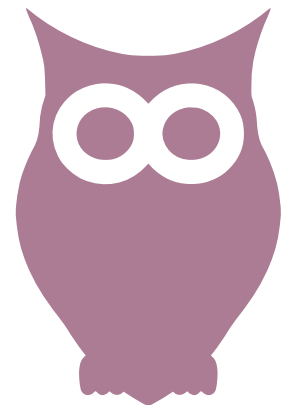  the sun's rays, he (Memnon) addressed him as well as he could.

Transl. Patricia A. Rosenmeyer

# Introducing GLAUx

Alek Keersmaekers

University of Leuven

# The GLAUx corpus

- The **G**reek **L**anguage **Au**tomated
- 26M tokens, literary Ancient Greek (with the papyri: 32M)
- Ranging from Homer to the Bible to mathematical texts to narrative texts to …
- 8th century BC – 3rd/4th century AD
- Automatic annotation:
  - Lemmas
  - Morphology
  - Syntax
  - Future: Semantics
- https://glaux.be/, https://github.com/alekkeersmaekers/glaux

# GLAUx annotation

- 'Trained' on manually annotated treebanks: AGDT, PROIEL, Gorman, Pedalion, Harrington, Yordanova (1,5M words) → about 5% of GLAUx
- Machine learning (AI) to predict annotations automatically
- Lemmas: ~99% accurate
- Morphology: ~97% accurate
- Syntax: ~80% accurate