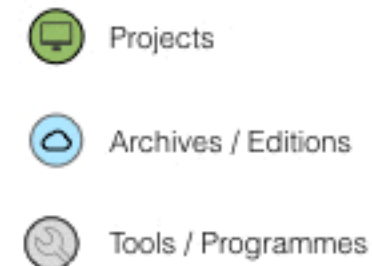
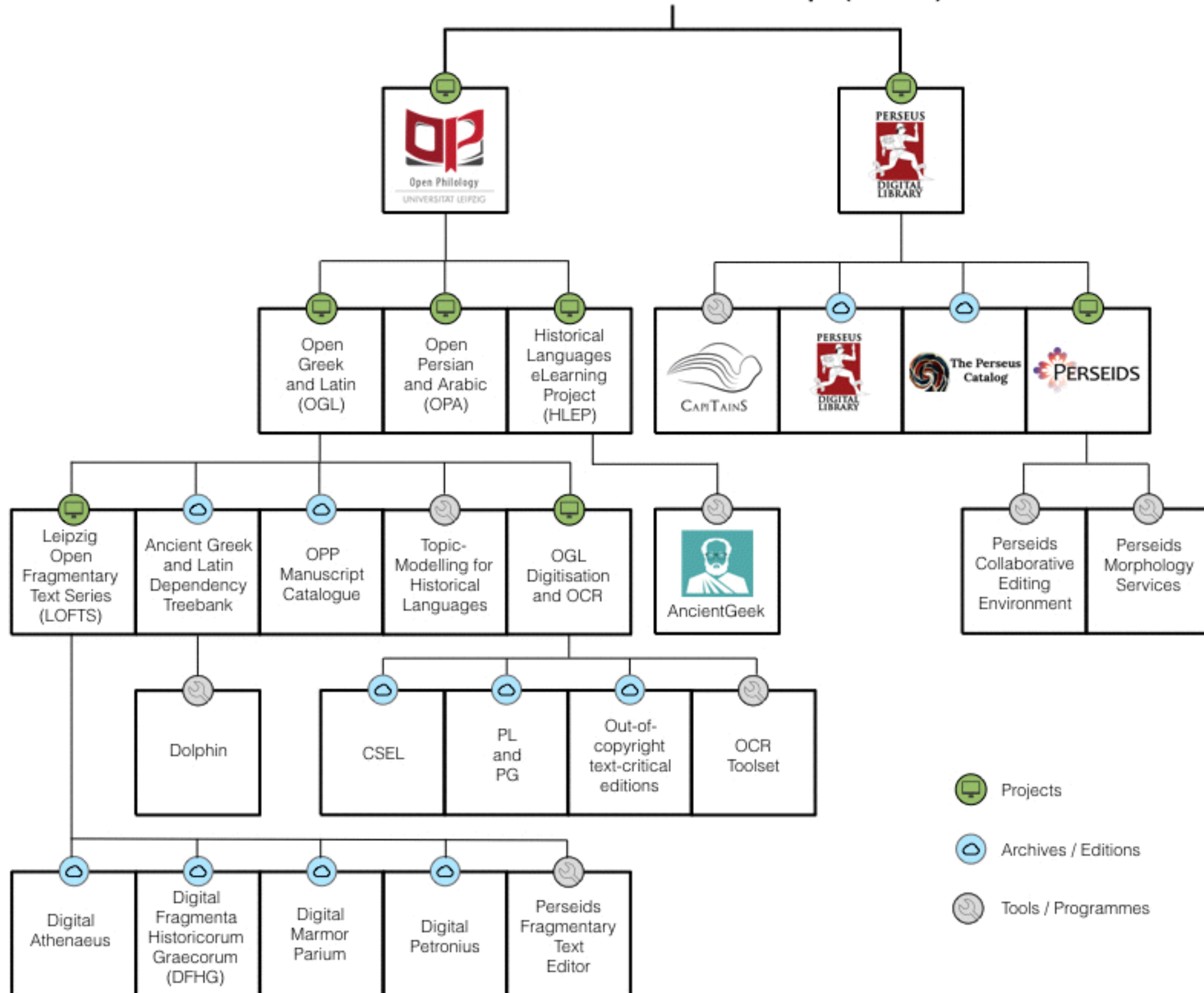


# Natural Language Processing of Historical Languages

Dr. Thomas Köntges,  
University of Leipzig,  
Open Philology Project,  
Perseus Research Group

# Perseus Research Group (PRG)



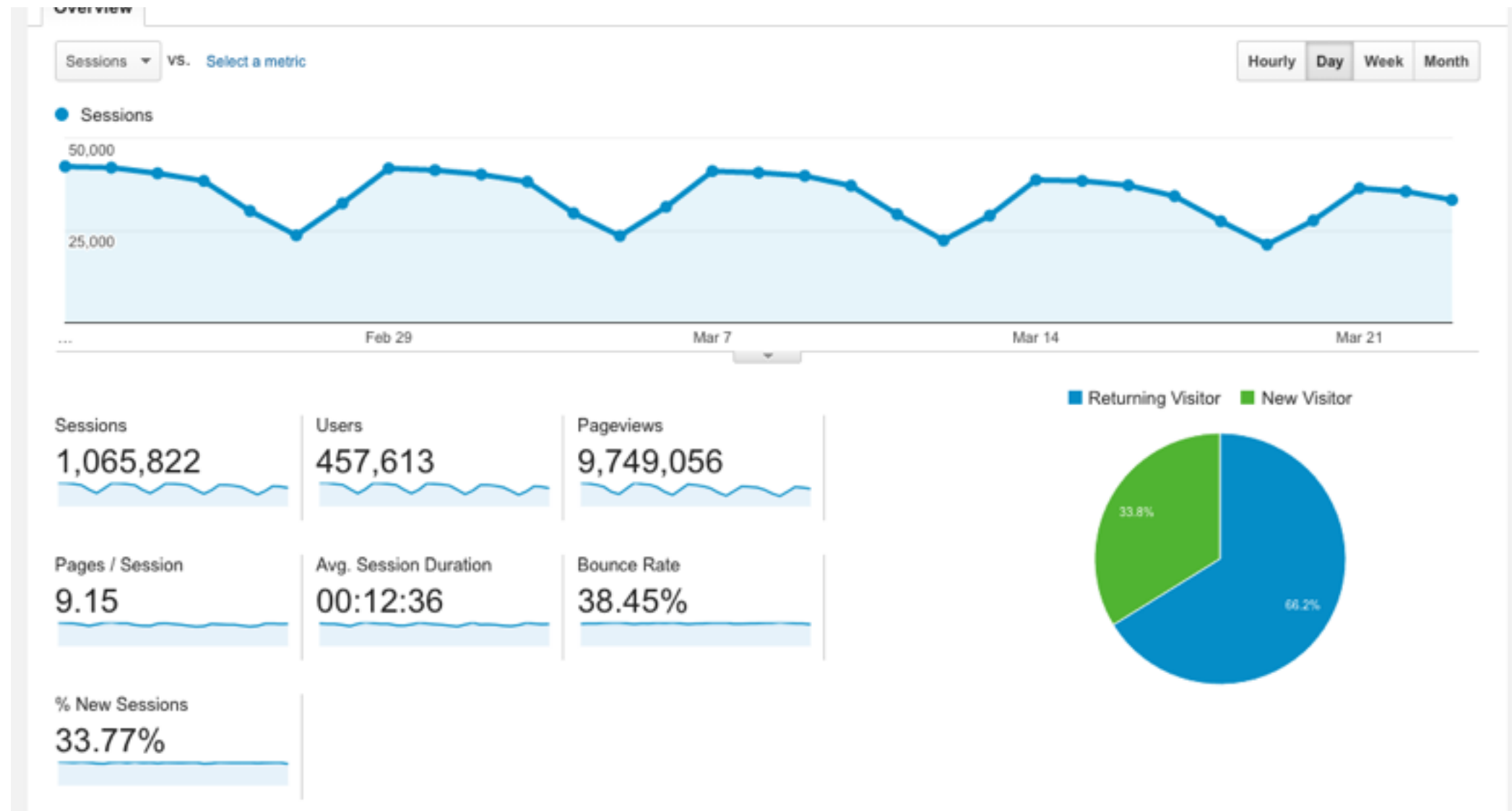
# Digital Renaissance

- Carolingian renaissance
- Italian renaissance
- digital renaissance
- access to literature, education, and reuse of knowledge





# Perseus Project visitors



# Natural Language Processing

# NLP

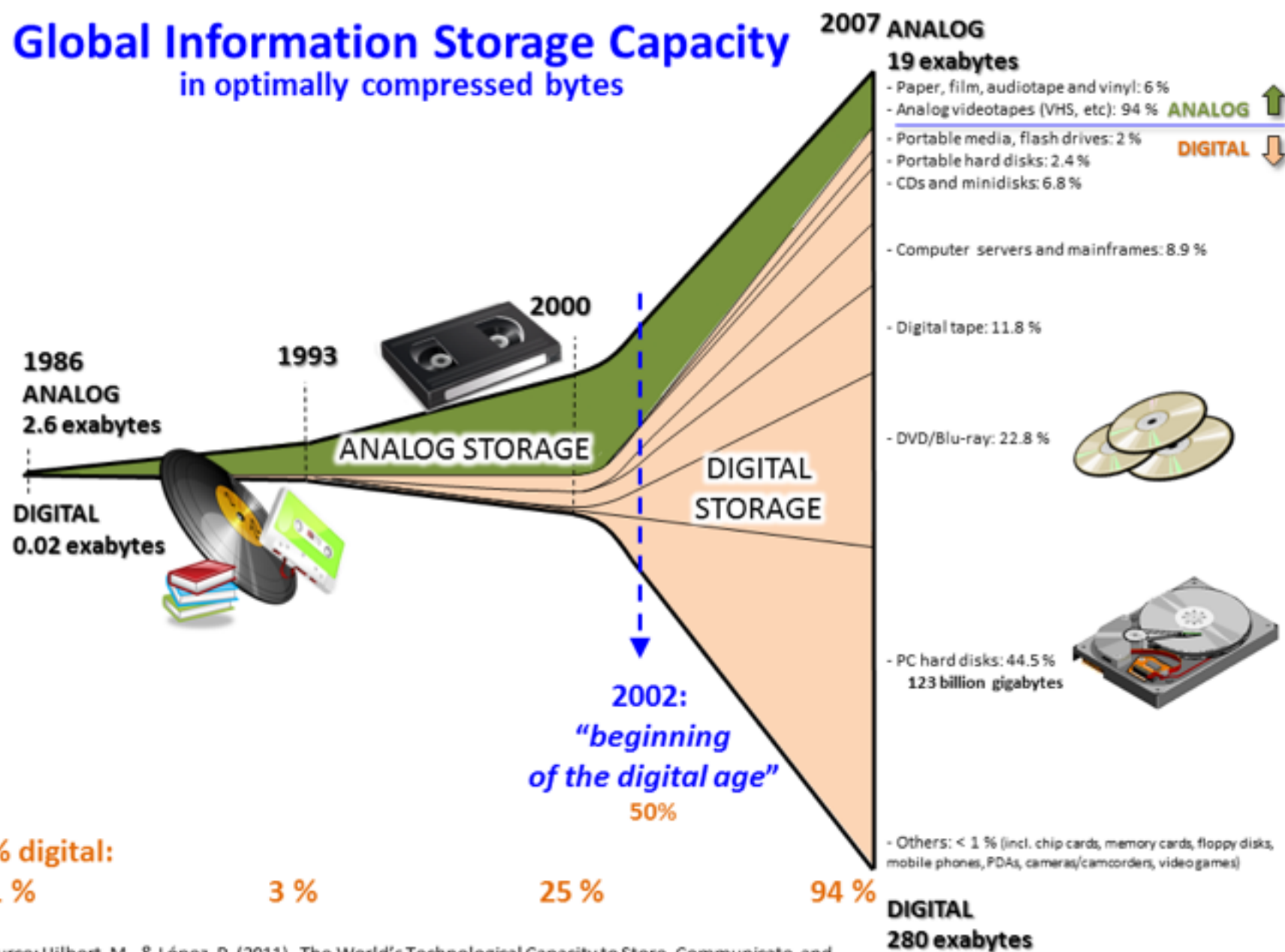
- NLP refers to the understanding and creation of human languages by a machine
- Today we focus on the former:
  - OCR
  - Machine Translation
  - Part-of-Speech Tagging
  - Topic and word sense disambiguation
  - Information Retrieval



BUT WHY??



# The Digital Age



Source: Hilbert, M., & López, P. (2011). The World's Technological Capacity to Store, Communicate, and Compute Information. *Science*, 332(6025), 60–65. <http://www.martinhilbert.net/WorldInfoCapacity.html>

# There is a lot of Data

The **2.8 Zettabytes** of data created in **2012** alone  
are the storage equivalent of  
**550 editions of Homer per day**  
**for each human being on the planet...**

# Textual Analysis

# Analysis

Split the texts into basic units.  
E.g. morphemes, words.



# Synthesis

Re-assemble those units  
into a new text (in the widest possible sense).

# Examples

- Metrical analysis
- Morphological analysis
- Tree-banking or other morpho-syntactic analyses
- Translation alignment
- Named-Entity mark-up

# Historical Languages

WHAT IS SPECIAL  
ABOUT THEM?



The complexity of a  
language changes over  
time.

# Grammar

- Indo-European languages have the tendency to become less morphologically complex.
- English barely 3 cases, German 4, Ancient Greek 4–5, Latin 5–6, Proto-Indo-European 8.
- English, German, Latin 2 Numbers, Ancient Greek 3.
- NB: Morphological complexity often replaced by syntactic or idiomatic complexity.
- NB: Languages often flatten if number of speakers is increased.

# PROBLEM OF DISTANCE & CONTEXT

Time flies like an arrow.

Fruit flies like an apple.

# Written Sources

A mono-directional communication of text-data, graphics, and meta-data from the past to the present.



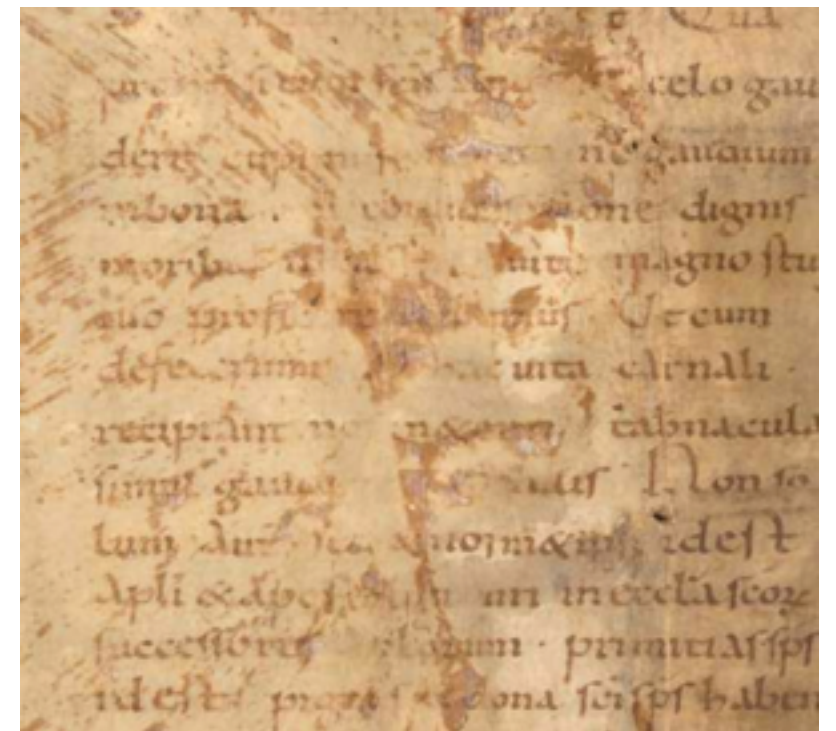


## TEXTUAL TRANSMISSION

A mono-directional communication  
of text-data, graphics, and meta-data  
from the past to the present.



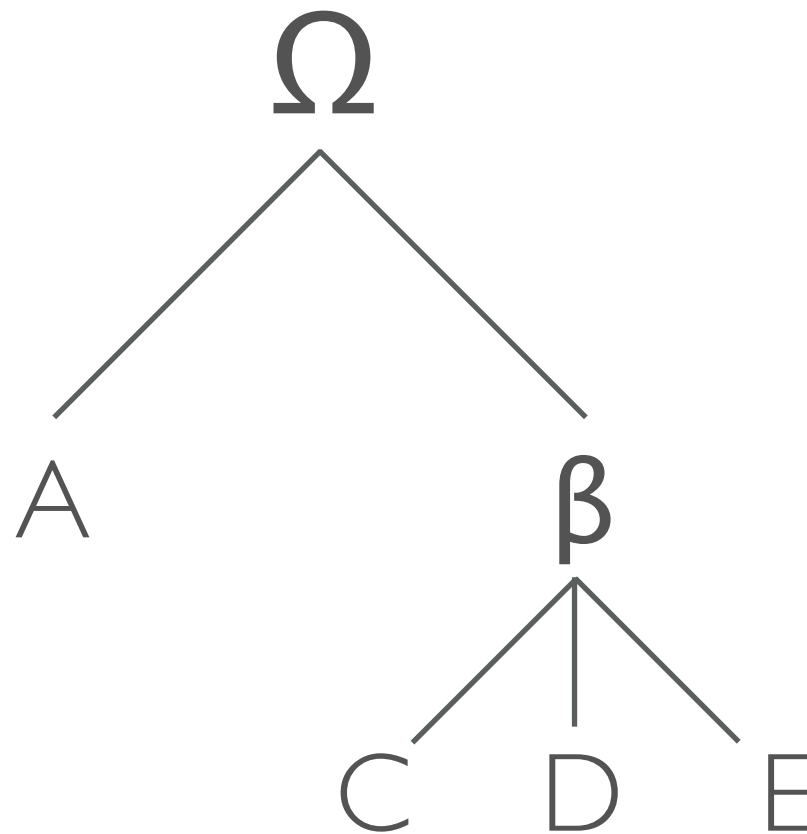
# Distribution is Essential



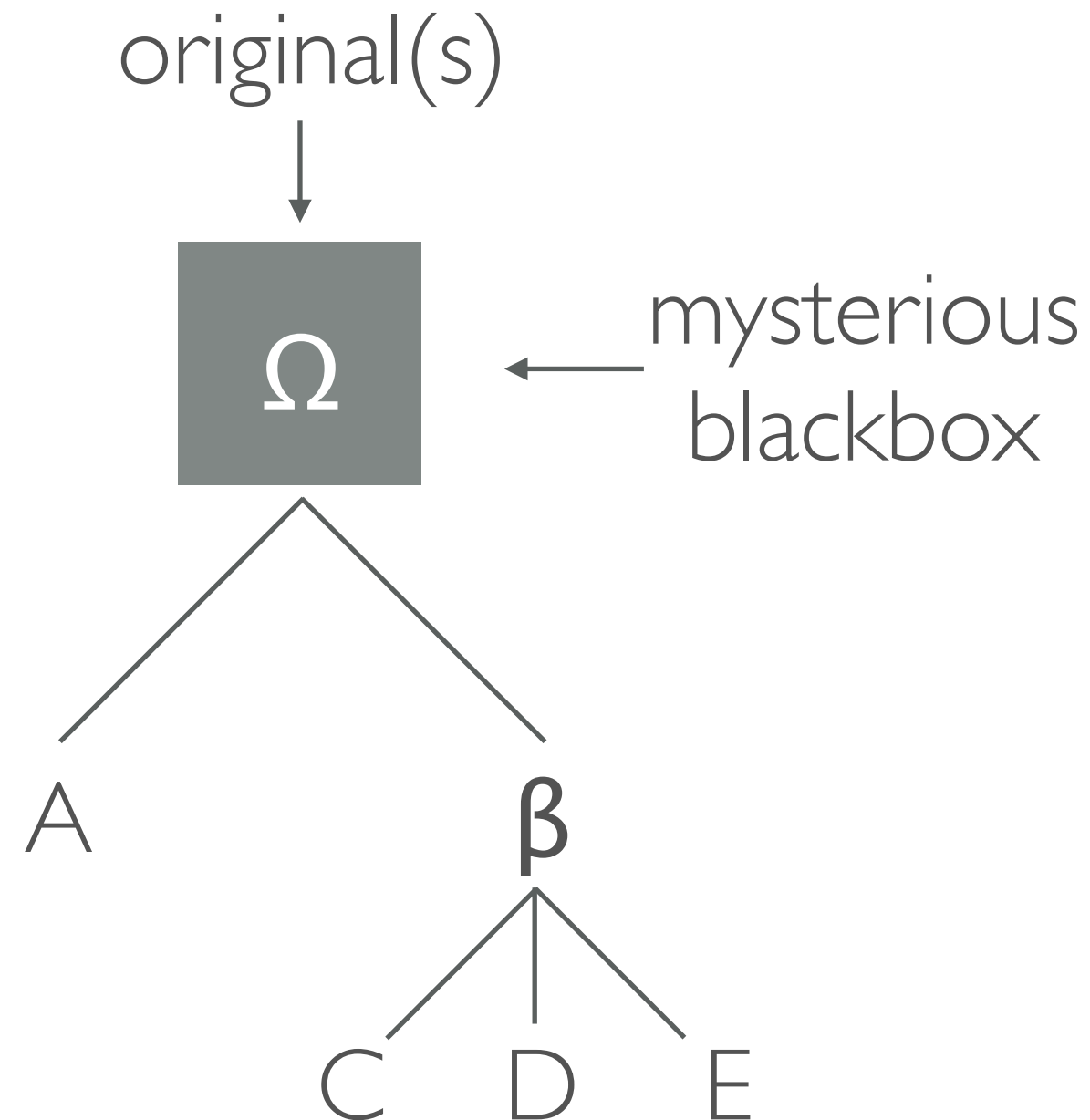
# Traditional Text-Criticism

- text-criticism is the skill of finding mistakes in ancient sources and the art of correcting — A. E. Housman (1922)
- the main task of text-criticism is to produce a text that resembles as closely as possible the autograph (original text). this process is called *constitutio textus* — P. Maas (1950)

# Archetype



# Archetype



# Archetype

an archetype represents a **virtual entity** that is the result and **summary** of all **unknown transmission processes** and changes the **original text** were subject to and that usually predates the oldest known/extant manuscripts of that text.

because an original text can have multiple originals (e.g. multitext, growing texts, canonical texts), it can have more than one archetype (as well as possible transmission-dependent hyparchetypi of each transmission line).

# Example: Petronius

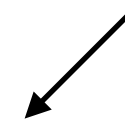
# Why Petronius?

- fragments of a fictional **prosimetric** text (35,000 words)
- fragments which look **connected** and build the base for our modern text (33,000)
- **references** in the text to **earlier passages** (23 references)
- **unconnected fragments** and carmina (51 fragments, at least 5 to 8 refer to pre-cena plot)
- transmitted along complex transmission lines in at least 60 manuscripts (4 main families) and hundreds of editions.



# Prosimetric

**prose**

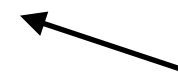


(§108.13) *data ergo acceptaque ex more patrio fide  
protendit ramum oleae a tutela navigii raptum, atque  
in colloquium venire ausa*

*‘Quis furor’ exclamat ‘pacem convertit in arma’*



**prose sentence's predicate**

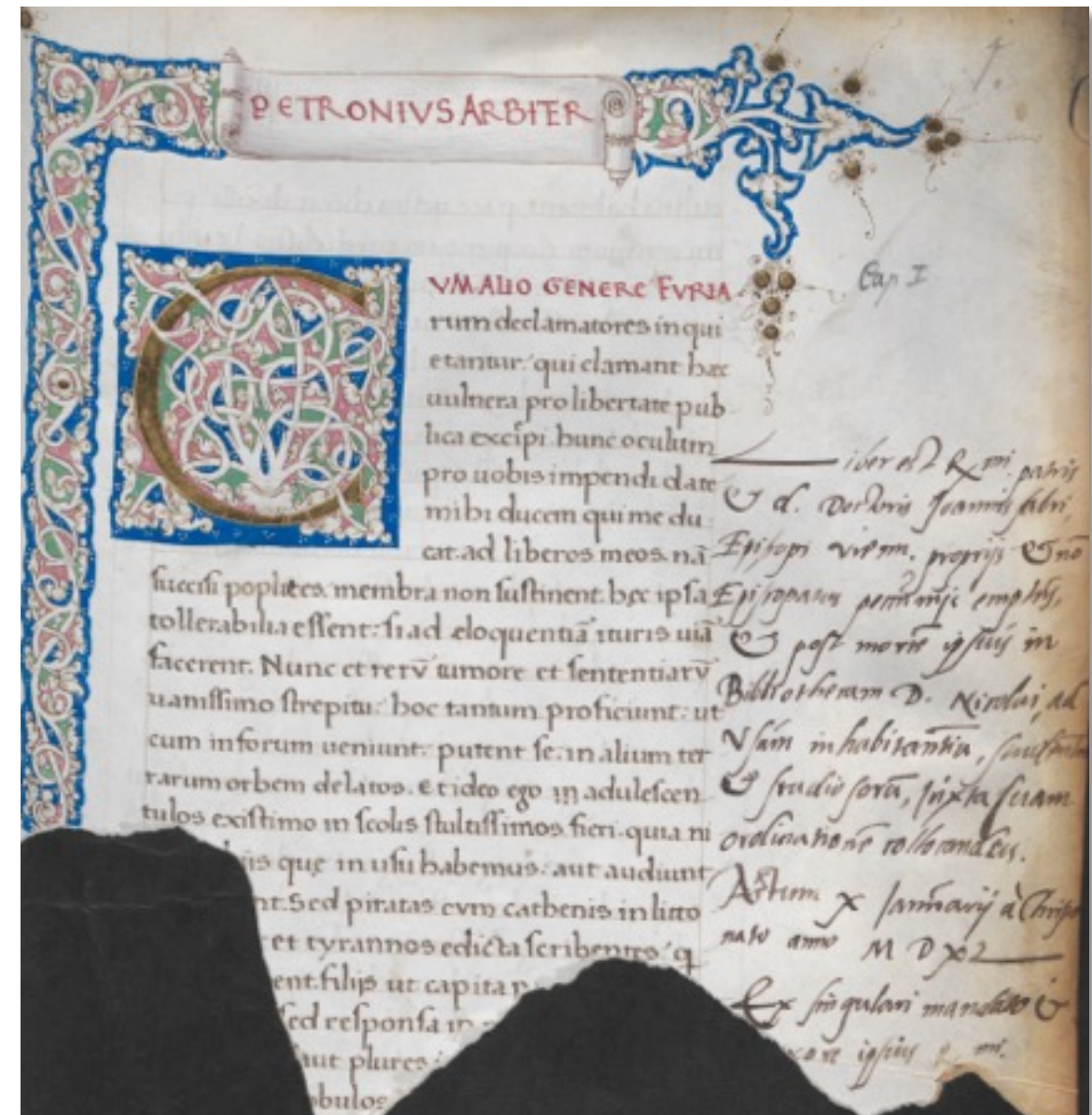
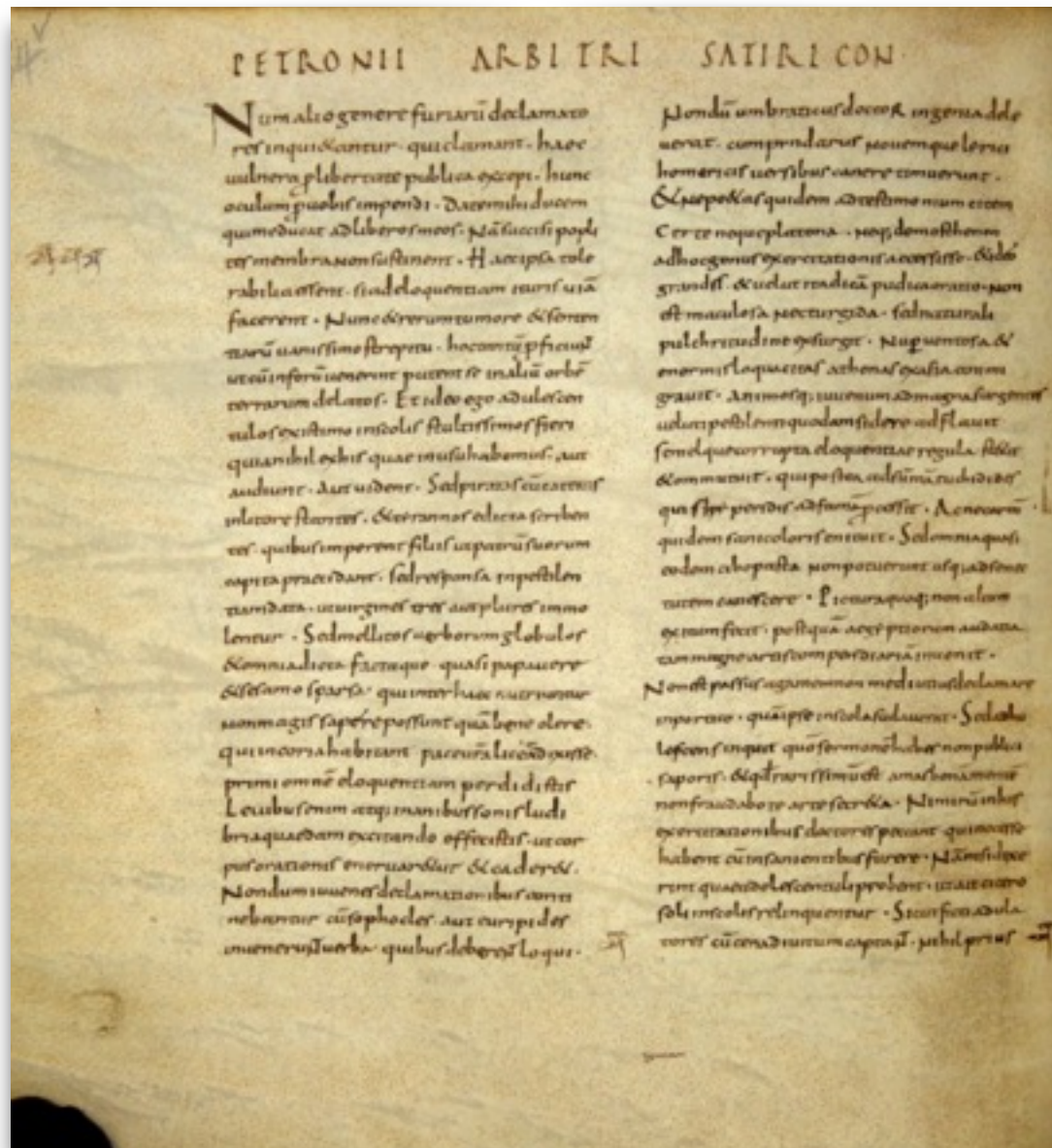


**verse**

# Connected Fragments

Text transmitted by the manuscripts					
Content	Chapters	Words	O	L	H
School Episode	1–5	700	yes	yes	
Brothel Episode	6–8	300	yes	yes	
Rape of Giton	9–9.5	100	only 9.5	yes	
Quarrel	9.6–11	300		yes	
Market Episode	12–15	700	poem 14.2	yes	
Quartilla's Entrance	16–19.1	500	yes	yes	
Orgy	19.2–26.6	1100	fr.	yes	
Bridge to <i>Cena</i>	26.7–27.1A	100			yes
Beginning <i>Cena</i>	27.1B–37.5	1800		yes *	yes
<i>Cena</i>	37.6–78	10500	only 55.4–6	fr.	yes
Ascyltos takes Giton	79–80	500	poem 80.9	yes	
Encolpius' Solitude	81–82	350	only 81.1–2	yes	
Meeting Eumolpus	83–84	350	fr.	yes	
Boy of Pergamon	85–87	600		yes	
Eumolpus about Literature	88	300	yes	yes	
<i>Halosis Troiae</i>	89	450		yes	
Reaction of Pedestrians	90	150		yes	
Reconnection with Giton	91–94	300	fr.	yes	
Eumolpus vs. <i>deversitor</i>	95–96	300	yes	yes	
Ascyltos' Search	97–99	650		yes	

# Manuscripts





# Editions

SATYRICON. 33  
dentibus oleum infuderat; & pueri, deterfis paulif-  
per oculis, redierant ad ministerium, quum <sup>93</sup> intrans  
cym-

mus inlustre, primum Archi-  
triclino Vinum adlatum, non vero  
ut biberet, sed ut Gullaret tan-  
tum: Vt autem Gullaret Archi-  
triclino aquam vinum fallam. ait D.  
Ioannes c. 11. Ubi indubie ejus in-  
dicatum fuisse officium animad-  
vertent Viri Docti, Indeque jam &  
Inscriptiones Lapidum clarescit:

A. LAGUNA. TRICLINARCHIA,  
Gonfalius.  
93. Intrans Cymbalistria, & con-  
crepans ara. Mulieres, quæ cytha-  
la quæiebant, pro analogia voca-  
mus CYMBALISTRIAS, vidi-  
musque eas sepius Romæ in scal-  
pturis Bacchanalium sacrorum, &  
quibus hæc damus:



Hinc autem adparet quomodo  
cymbala colliderentur, quod in-  
telligens Amalarius Fortunatus  
scripsit lib. III. c. III. cymbala  
invicem tanguntur ut sonent  
ideo à quibusdam labijs nostris com-  
parata sunt. Dioscorides lib. IV.  
testatur cymbala ejus figuræ fuisse,  
qua

Varia lectiones ex collatione v. c.

Pagina 1. linea 19. & tyrannos p. 2. l. 5. quibus  
loqui debemus l. 8. deest certe l. 14. emigravit,  
l. 20. vsque ad sen. p. 3. l. 3. quam id quod l. 5., ni  
si quasdam insidias auribus fecerint: sic eloquen-  
tiæ magistri, tanquam piscator, qui nisi eam l.  
12. impellunt, l. 15. vt stud. l. 28. prius meræ l.  
29. palleat p. 4. l. 3. redimitus histrioni l. 4. crini

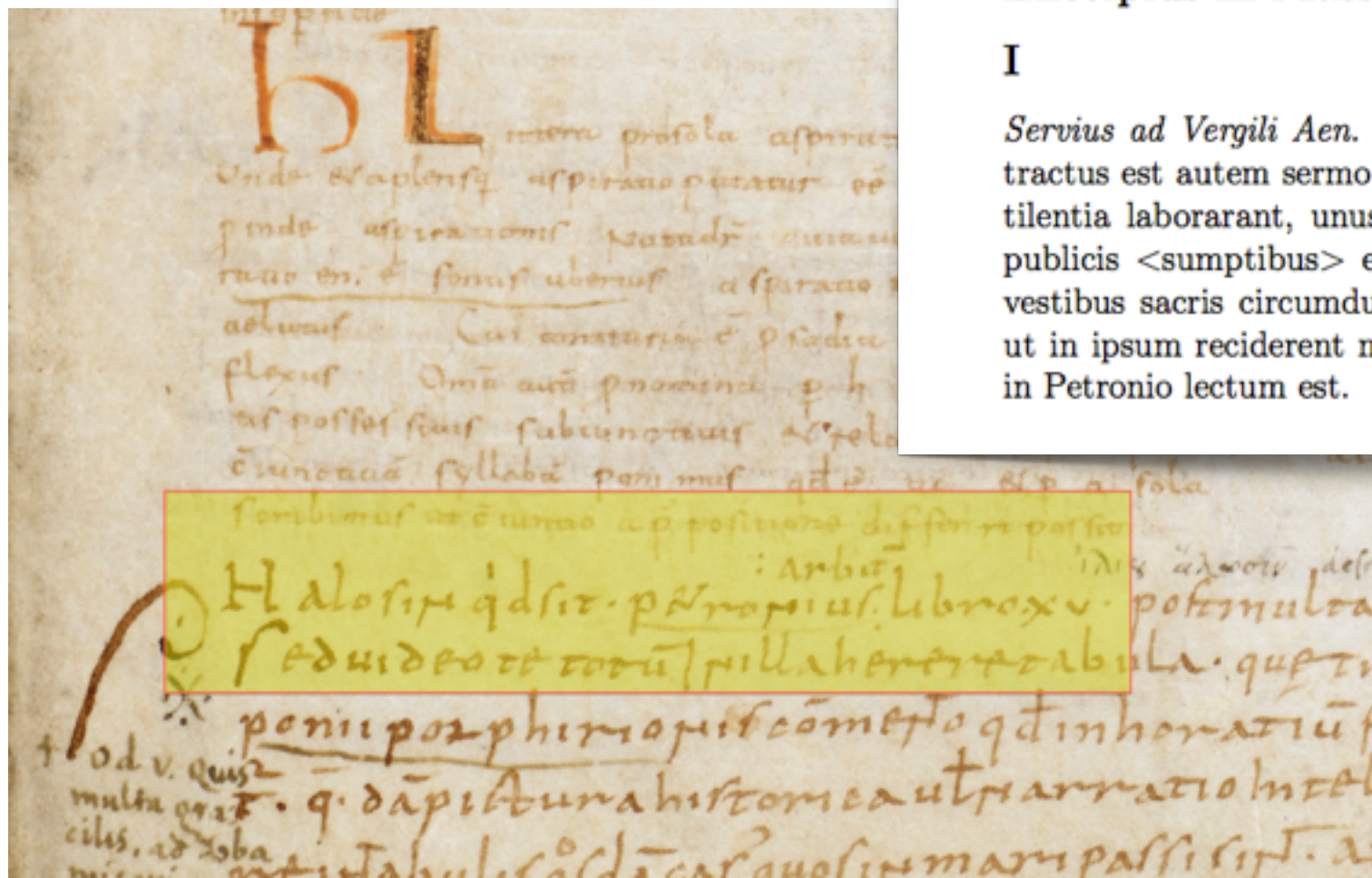
# References

## FRAGMENTA PETRONIANA

### Encolpius in Massilia

#### I

*Servius ad Vergili Aen. III 57: auri sacra fames*] sacra id est execrabilis. tractus est autem sermo ex more Gallorum. nam Massilienses quotiens pestilentia laborarant, unus se ex pauperibus offerebat, alendus anno integro publicis <sumptibus> et purioribus cibis. hic postea ornatus verbenis et vestibus sacris circumducebatur per totam civitatem cum execrationibus, ut in ipsum reciderent mala totius civitatis, et sic proiciebatur. hoc autem in Petronio lectum est. 5





# Self-References

## FRAGMENTA DE PETRONIO APPROBATA

### Conspectus (FrPA I)

#### FrPA I (= 81.3b–5)

- 81.3b | effugi iudicium, harenae imposui, hospitem occidi, ut inter <tot> audaciae *L*  
nomina mendicus, exul, in deversorio Graecae urbis iacerem desertus?...
- 4 adulescens omni libidine impurus et sua quoque confessione dignus exilio,  
stupro liber, stupro ingenuus, cuius anni ad tesseram venierunt, quem tam-
- 5 quam puellam conduxit etiam qui virum putavit. quid ille alter? qui *5*  
[tamquam] die togae virilis stolam sumpsit, qui ne vir esset a matre per-  
suasus est, qui opus muliebre in ergastulo fecit, qui postquam conturbavit  
et libidinis suae solum vertit, reliquit veteris amicitiae nomen et, pro pudor,  
tamquam mulier secutuleia unius noctis tactu omnia vendidit.

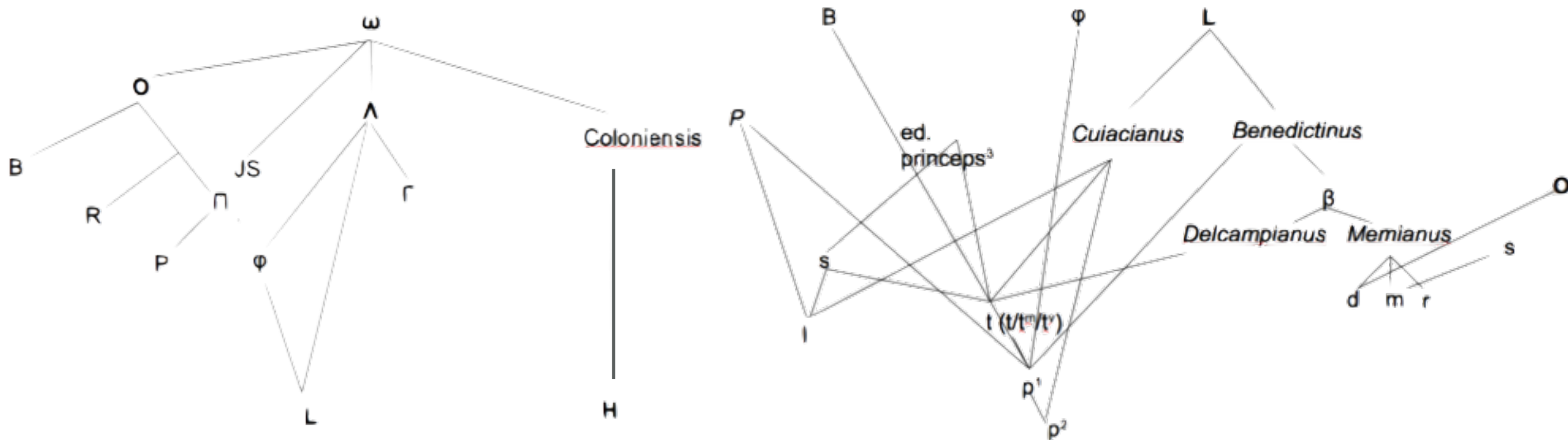
### The sanctuary of Priapus (FrPA XXI–XXII)

#### FrPA XXI (= 16.3)

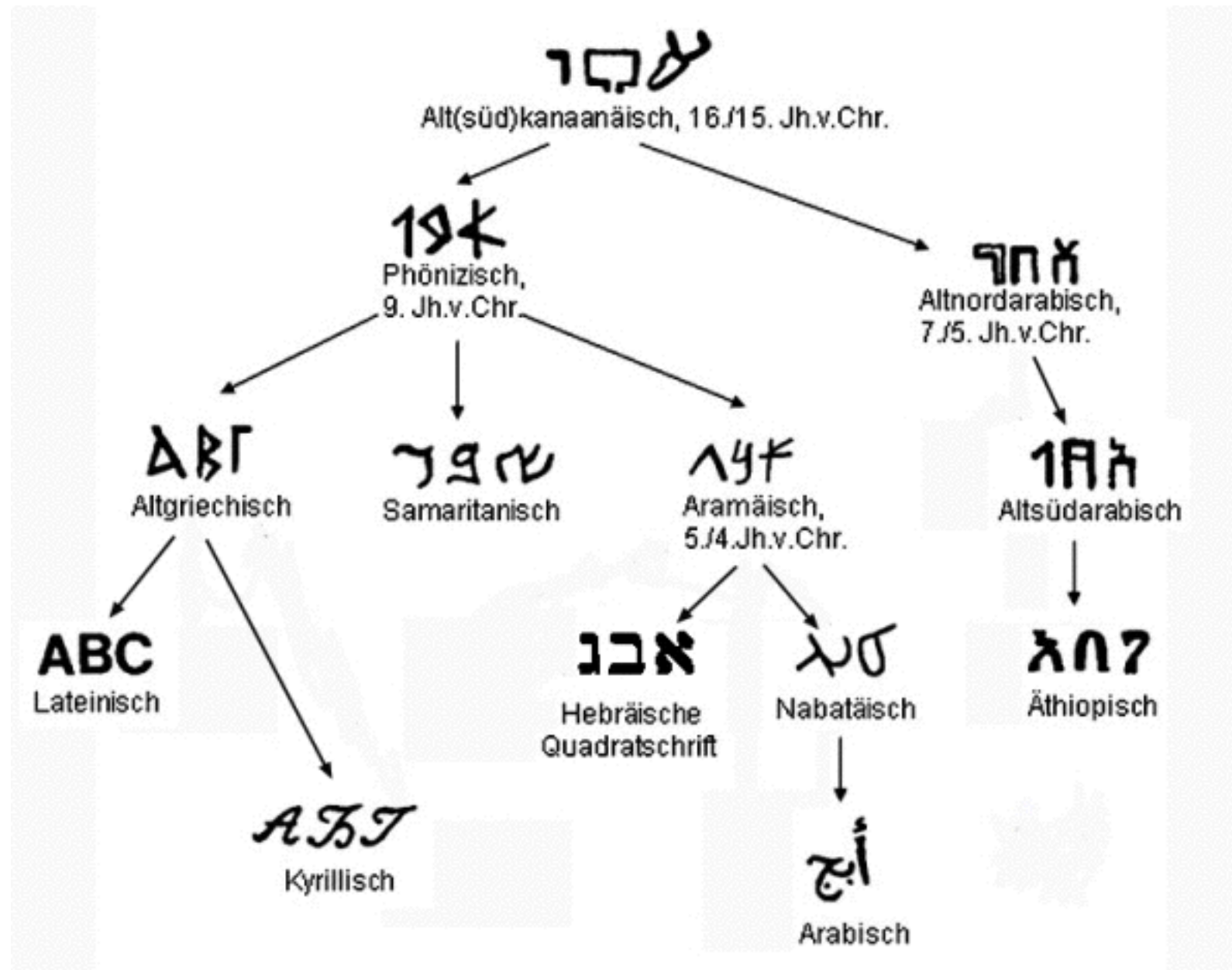
- 16.3 | 'ego sum ancilla Quartillae, cuius vos sacrum ante cryptam turbastis. *LO*

# Reconstructions

- Bücheler 1862, editio maior (14 manuscripts)
- Müller 2009, Satyricon reliquiae (24 manuscripts)



# Code Change





# Code Change

Trajan	A B C D E F G H I K L M N O P Q R S T V Y Z
Rustic	A B C D E F G H I K L M N O P Q R S T V X Y
Greek Uncial	Α Β Γ Δ Ε Ζ Η Θ Ι Κ Λ Μ Ν Ξ Ο Π Ρ Σ Τ Υ Φ Χ Ψ Ω
Uncial	A B C D E F G H I K L M N O P Q R S T U X Y
Half-Uncial	a b c d e f g h i k l m n o p q r s t u x y
Visgothic	α β γ δ ε ζ η θ ι κ λ μ ν ο ρ ρ σ τ υ χ ψ ζ
Luxeuil	u b c d e f g h i k l m n o p q r r s t u x y
Beneventan	α β γ δ ε ζ η ι κ λ μ ν ο ρ q r s t u x y z
Caroline	A B C D E F G H I J K L M N O P Q R S T U X Y
Insular	α β γ δ ε ζ η ι κ λ μ ν ο ρ q r s t u v r x y z
Protogothic	a b c d e f g h i j k l m n o p q r s t u v w x y z
Textualis quadrata	a b c d e f g h i j k l m n o p q r s t u v w x y z
Fraktur	a b c d e f g h i j k l m n o p q r s t u v w x y z
Humanist	a b c d e f g h i j k l m n o p q r s t u v w x y z
Times	a b c d e f g h i j k l m n o p q r s t u v w y z

# Examples

# Combination of Analyses

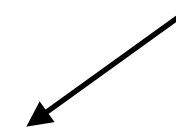
# CTS/CITE Architecture

# canonical text services

- stable identifier for canonical texts
- traditional citation  $\neq$  canonical citation
- sentence level  $\neq$  canonical citation (verse, prose, prosimetric texts, religious texts...)
- what the digamma, is a canonical citation then?

# CTS/CITE URNs

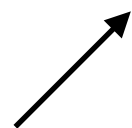
**CTS URN**



Namespace	Specific	CTS Namespace	TextGroup	Work	Version	Exemplar	Citation	Subreference
urn	:cts	:greekLit	:tlg0012	.tlg001	.msA	.thosJeff	:1.2	@ούλομέν

urn:cite:CITENAMEPSACE:COLLECTION.OBJECTID

**CITE URN**



urn:cite:hmt:msA.12r

# Homer multitext project

# HMT

Digital services from the Homer Multitext project

[HMT: home](#) | [HMT Digital: home](#) | [browse MSS](#) | [scholia](#)

## HMT Digital: Scholia Reader

8 Scholia for 1.1

Manuscript [urn:cts:greekLit:tlg0012.tlg001.msA:1.1](#)

Homeric epic *Iliad* A

§1

1 Μῆνιν ἄειδε θεὰ Πηληϊάδεω Ἀχιλῆος

( [urn:cts:greekLit:tlg0012.tlg001.msA:1.1](#) )

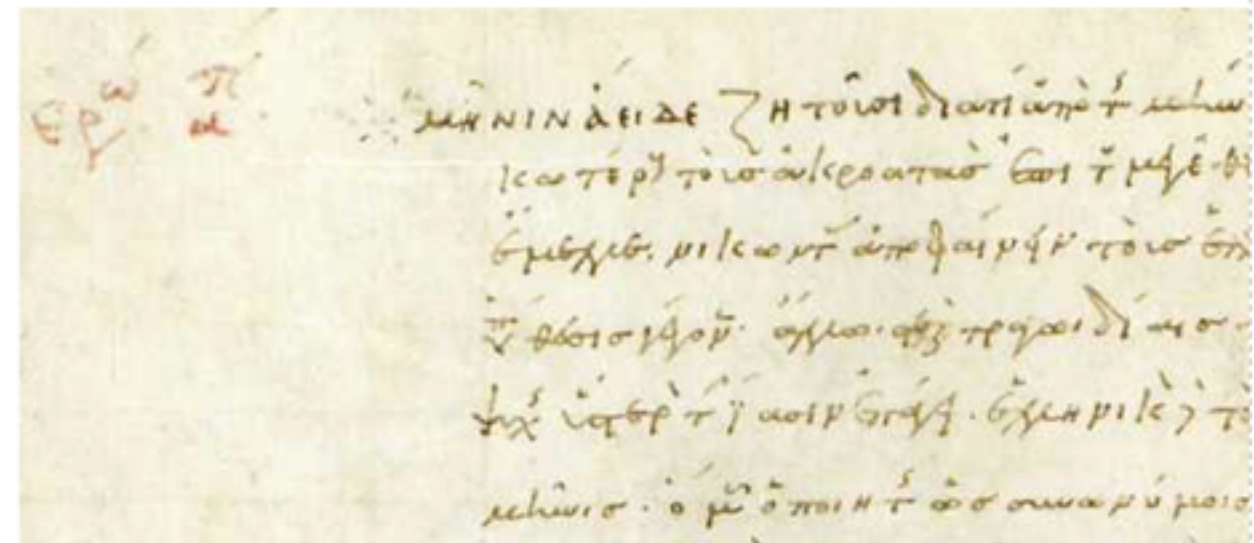
Scholia to the Iliad *Main marginal scholia* msA

§1

§1

μῆνιν ἄειδε

ζητοῦσι δια τί ἀπὸ τῆς μῆνιδος ἤρξατο οὕτως δυσφήμου ὀνόματος·  
δια δύο ταῦτα· πρῶτον μὲν ἵν' ἐκ τοῦ πάθους ἀπκατάρρεῦση τὸ  
τοιούτο μόριον τῆς ψυχῆς καὶ προσεκτικωτέρους τοὺς ἀκροατὰς ἐπὶ  
τοῦ μεγέθους ποιήσῃ καὶ προσεθίζῃ φέρειν γενναίως ἡμᾶς τὰ πάθη.  
μέλλων πολέμους ἀπαγγέλλειν· δεύτερον, ἵνα τὰ ἐγκώμια τῶν  
Ἑλλήνων πιθανώτερα ποιήσῃ· ἐπεὶ δὲ ἔμελλε, νικωντας ἀποφαίνειν  
τοὺς Ἑλληνας, εἰκότως οὐ κατατρέχει ἀξιοπιστότερον ἐκ τοῦ μὴ  
πάντα χαρίζεσθαι τῷ ἐχθρῶν ἐπαίνῳ· ἤρξατο μὲν ἀπὸ μῆνιδος





# Manage Citations



The Perseus Catalog interface displays search results for the query "tg\_no\_token:[E TO F]". The search bar shows the query and a "Search" button. The results section lists 17 items, with the first item being "Catoptrica" by Euclid. The interface includes a "Browse by" sidebar with categories like Author, Work Title, and Work Original Language. The footer contains logos for Tufts University and The University of Leipzig, along with links for "About the Catalog" and "Contact Us".

**The Perseus Catalog**

Search History  
Author List  
Help

All Fields  
tg\_no\_token:[E TO F] Search

You searched for: tg\_no\_token:[E TO F] x Start Over

Author > Euclid. x

1 - 17 of 17 20 per page Sort by relevance

**1. Catoptrica**

URN: urn:cts:greekLit:tlg1799.tlg011  
Author: Euclid  
Language: Greek, Ancient (to 1453)

**2. Data**

Perseus Digital Library  
Tufts University Medford, MA, USA The University of Leipzig Leipzig, Germany

About the Catalog Contact Us




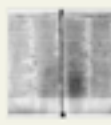



The Perseus Catalog, Version 1.0

# Manage the Content

DH-Leipzig Welcome, Thomas Koentges [Log Out](#)

**Browse Items (19 total)**

[Add an Item](#) [Show Details](#) [Search Items](#) [Edit](#) [Delete](#) [Quick Filter](#) 1 of 2 [>](#)

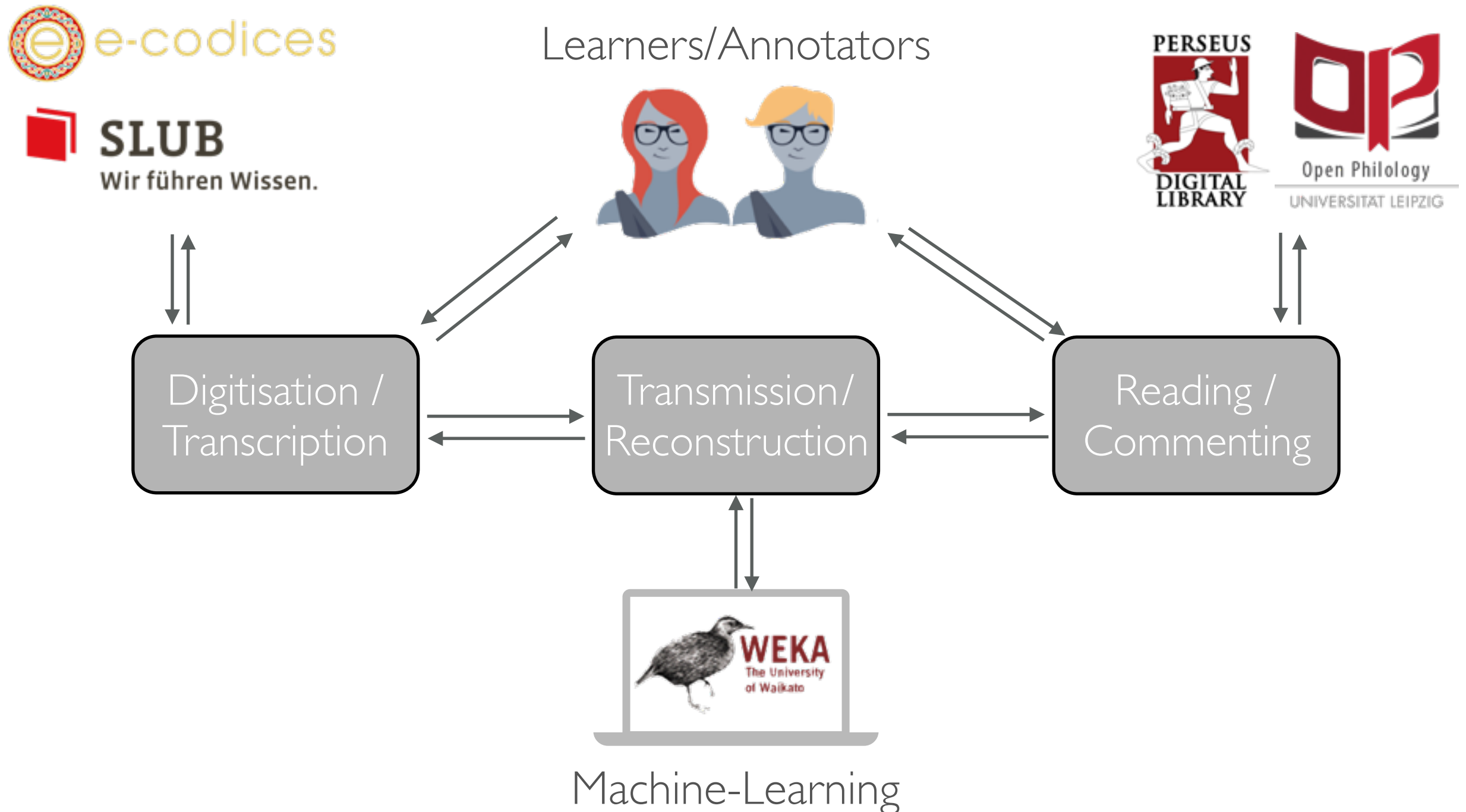
<input type="checkbox"/>	Title	Creator	Type	Date Added
<input type="checkbox"/>	 <b>Biblioteca Estense, alfa.t.8.21 (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> No document identifier.			Jun 23, 2015
<input type="checkbox"/>	 <b>Bayerische Staatsbibliothek, Cod. graec. 427 (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> No document identifier.			Jun 23, 2015
<input type="checkbox"/>	 <b>BNF, Lat 7975, 81v - 82r (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> urn%3Acite%3Aogl%3Aabnf_lat7975			Jun 22, 2015
<input type="checkbox"/>	 <b>BNF, Lat 7647 110v-112v (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> urn%3Acite%3Aogl%3Aabnf_lat7647		Still Image	Jun 22, 2015
<input type="checkbox"/>	 <b>Plut. 47.31, f.128r-150v (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> urn%3Acite%3Aogl%3Aplut47_31			Jun 21, 2015
<input type="checkbox"/>	 <b>Plut. 37.25 (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> urn%3Acite%3Aogl%3Aplu37_25			Jun 21, 2015
<input type="checkbox"/>	 <b>Leiden VLF 111 (Private)</b> <a href="#">Details</a> · <a href="#">Edit</a> · <a href="#">Delete</a> urn%3Acite%3Aogl%3Aleiden_vlf111			Jun 21, 2015



CTS/CITE architecture:  
find the logical structure of  
a text and make it citable.

# Textual Transmission and Machine Learning / Data Mining

# arXe.type Workflow



# Optical Character Recognition / Transcription



# OCR

booktif/0016/01000e.bin.png

**stularemus misit in faciem Ascylti tunicā, & libe**

stulatemus misit in faciemu Afelti tunicā, & libo

booktif/0016/01000f.bin.png

**ratos querela iussit pallium deponere quod fo-**

tatos querela iussit pallium deponete quod fo-

booktif/0016/010010.bin.png

**lum hitem faciebat, & recuperato ut putabamus**

lumu fitem ffaciebat, & recupetato lt putabamus

booktif/0016/010011.bin.png

**thesauro in diuersorium principes abimus, præ-**

thesauro tn diuetforium ptrcipites abimus, præ-

booktif/0016/010012.bin.png

**clusisq; foribus ridere acumē non minus concio**

clufisqp foribus ridere acumē non minus concio

# OCR

<sup>1</sup>  
 Ὁ μὲν οὖν Αὐγουστος ἀπεβίω ὡς εἴρηται, τὴν <sup>PI545</sup>  
 δὲ μοναρχίαν ὁ Τιβέριος διεδέξατο. ὃς εὐπατρίδης <sup>WH171</sup>  
 μὲν ἦν καὶ πεπαιδευτο, τὴν δὲ γνώμην ἦν ποικιλώ-  
 τatos, ἐναντίους τῇ προαιρέσει τοὺς λόγους ποιού-  
 5 μενος. ὧν γὰρ ἐβούλετο τάναντία ἔλεγεν, ἅλλα μὲν  
 κεύθων ἐνὶ φρεσίν, ἅλλα δὲ λέγων· καὶ ὀργίζεσθαι  
 προσποιούμενος ἐν οἷς οὐκ ὠργίζετο, καὶ ἐν οἷς ἐθυ-  
 μοῦτο σχηματιζόμενος ἐπιείκειαν· καὶ ὡς οἰκειότατον  
 ἑώρα τὸν ἔχθιστον, καὶ ὡς ἀλλοτριωτάτῳ προσεφέ-  
 10 ρετο τῷ φιλάτῳ. καὶ οὐκ ἤξλου τοῖς ἄλλοις δῆλον  
 εἶναί οἱ τὸ φρόνημα, τοῦτο προσήκειν τῷ αὐταρχοῦντι  
 φρονῶν. καὶ εἴτε τις ἠναντιοῦτο οἷς ἔλεγεν εἴτε μὴν D  
 καὶ συνήνει, μεμίσητο.

Τέως δ' οὖν εἰς τὰ στρατόπεδα καὶ εἰς τὰ ἔθνη  
 15 πάντα ὡς αὐτοκράτωρ ἀντίκα ἐπέστειλε, μὴ λέγων  
 αὐτοκράτωρ εἶναι· ψηφισθὲν γὰρ αὐτῷ καὶ τοῦτο  
 μετὰ τῶν ἄλλων ὀνομάτων οὐκ ἐδέξατο. καὶ τὰ τῆς  
 ἀρχῆς διοικῶν ἅπαντα, μηδὲν αὐτῆς δεῖσθαι ἔλεγε,  
 καὶ ταύτης ἐξίστασθαι ἐκομψεύετο καὶ διὰ τὴν ἡλι-  
 20 κίαν, ἔξ γὰρ καὶ πεντήκοντα ἔτων ἦν, καὶ δι' ἀμ-  
 βλυωπίαν· πλείστον γὰρ ἐν σκότει βλέπων ἐλάχιστα  
 τὴν ἡμέραν εἴωρα. εἶτα κοινωνοὺς ἦτει τῆς ἀρχῆς  
 καὶ συνάρχοντας, οὐδὲν τούτων ποιῆσαι μέλλον, ἀλλ' ΠΙ546

ZONARAS III.

1

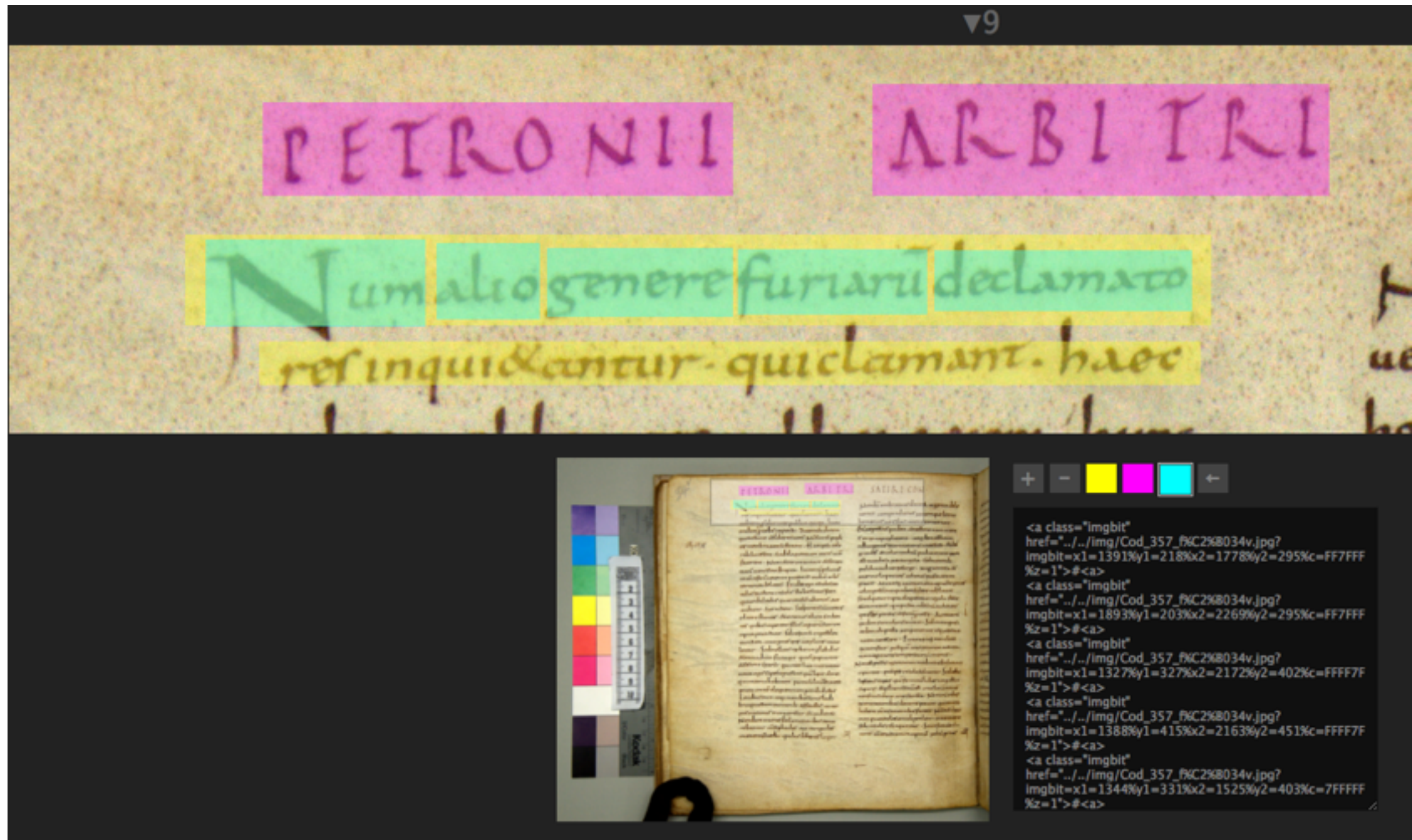
.PI515

Δ μὲν οὖν 7ύγουστος ἀπεβίω ὡς εἴρηται, τὴν,  
δὲ μοναρχίαν ὁ Ιβέριος διεδέξατο. ὃς εὐπατρίδης VW II171  
μὲν ἦν καὶ πεπαίδευτο, τὴν δὲ γνώμην ἦν ποικιλώ-  
τατος, ἐναντίους τῇ προαιρέσει τοὺς λόγους ποιού-  
5 μενος. ν γὰρ ἐβούλετο τάναντία ἔλεγεν, ἄλλα μὲν  
κεύθων ἐνὶ φρεσίν, ἄλλα δὲ λέγων· καὶ ὀργίζεσθαι  
προσποιοῦμενος ἐν οἷς οὐκ ὠργίζετο, καὶ ἐν οἷς ἐθυ-  
μούτο σχηματιζόμενος ἐπείκειαν· καὶ ὡς οἰκειότατον  
ἔώρα τὸν ἔχθιστον Pa, PPMέ ὡς ἀλλοτριωτάτῳ προσεφέ-  
10 ρετο τῷ φιλάτῳ. καὶ οὐα ἡξιου τοῖς ἄλλοις δῆλον  
εἶναί οἱ τὸ φρόνημα, τοῦτο προσήκειν τῷ αὐταρχοῦντι  
φρονῶν. καὶ εἴτε τις ἠναντιοῦτο οἷς ἔλεγεν εἵτε μὴν I)  
καὶ συνήνει, μεμίοητο.

ἕως δ' οὖν εἰς τὰ στρατόπεδα καὶ εἰς τὰ ἔθνη  
15 πάντα ὡς αὐτοκράτωρ αὐτίκα ἐπέστειλε, μὴ λέγων  
αὐτοκράτωρ εἶναι· ψηφισθὲν γὰρ αὐτῷ καὶ τοῦτο  
μετὰ τῶν ἄλλων ὀνομάτων οὐκ ἐδέξατο. καὶ τὰ τῆς  
ἀρχῆς διοικῶν ἅπαντα, μηδὲν αὐτῆς δεῖσθαι ἔλεγε,  
καὶ ταύτης ἐξίστασθαι ἐκομψεύετο καὶ διὰ τὴν ἡλι-  
20 κίαν, ἧν γὰρ καὶ πεντήκοντα ἐτῶν ἦν, καὶ δι' ἀμ-  
βλυωπίαν· πλειστον γὰρ ἐν σκότει βλέπων ἐλάχιστα  
τὴν ἡμέραν ἑώρα. εἶτα κοινωνοὺς ἦtet τῆς ἀρχῆς  
καὶ συνάρχοντας, οὐδὲν τούτων ποιῆσαι μέλλων, ἀλλ' PI516  
Cap. 1. Dionis Historiae Romanae 1. 57, c. 1-13.  
eoNARA8 1II.



# Transcription



# Transcription

book/0001/010001.bin.png

Alpibus excludo uincendo certior exul ..

Alpibus excludo uincendo certior exul

book/0001/010002.bin.png

Sanguine germano. Lxq triumphis sexaginta

Sanguine germano. Lxq triumphis

book/0001/010003.bin.png

Esse nocens coepi. q̃q̃ quos gloria terret

esse nocens coepi. q̃q̃ quos gloria terret

book/0001/010004.bin.png

Aut qui fūt qui bellā uident mercedib9 epte

Aut qui fūt qui bellā uident mercedib<sup>9</sup> epte

# Weka

- Weka 3.7.11: open source data mining software in Java
- a collection of machine-learning algorithms for data-mining tasks
- although Weka is a powerful tool, classifying the manuscripts will be a challenge

# Weka

## @relation PetroniusManuscripts

@attribute ID numeric

**@attribute Variant1 {Cum,Num,om.,Quum}**

@attribute Variant2 {Nec,Non,om.}

@attribute Variant3 {Essent,om.}

@attribute Variant4 {ituris,nuris,om.,turis}

@attribute Variant5 {facent,facerent,om.}

@attribute Variant6 {rerum,runci,verborum,om.}

@attribute Variant7 {et,om.}

@attribute Variant8 {strepita,strepitum,om.}

@attribute Variant9 {terrarium\_orbem,orbem\_terrarum,om.}

@attribute Variant10 {qui,quia,quorum,om.}

@attribute Variant11 {his,iis}

@attribute Variant12 {et,om.,sed}

@attribute Variant13 {et,om.,sed}

@attribute Variant14 {pestilentia,pestilentiam,om.}

@attribute Variant15 {omnem,omnium,om.}

@attribute Variant16 {quidem,om.,solum}

@attribute Class {L,O}

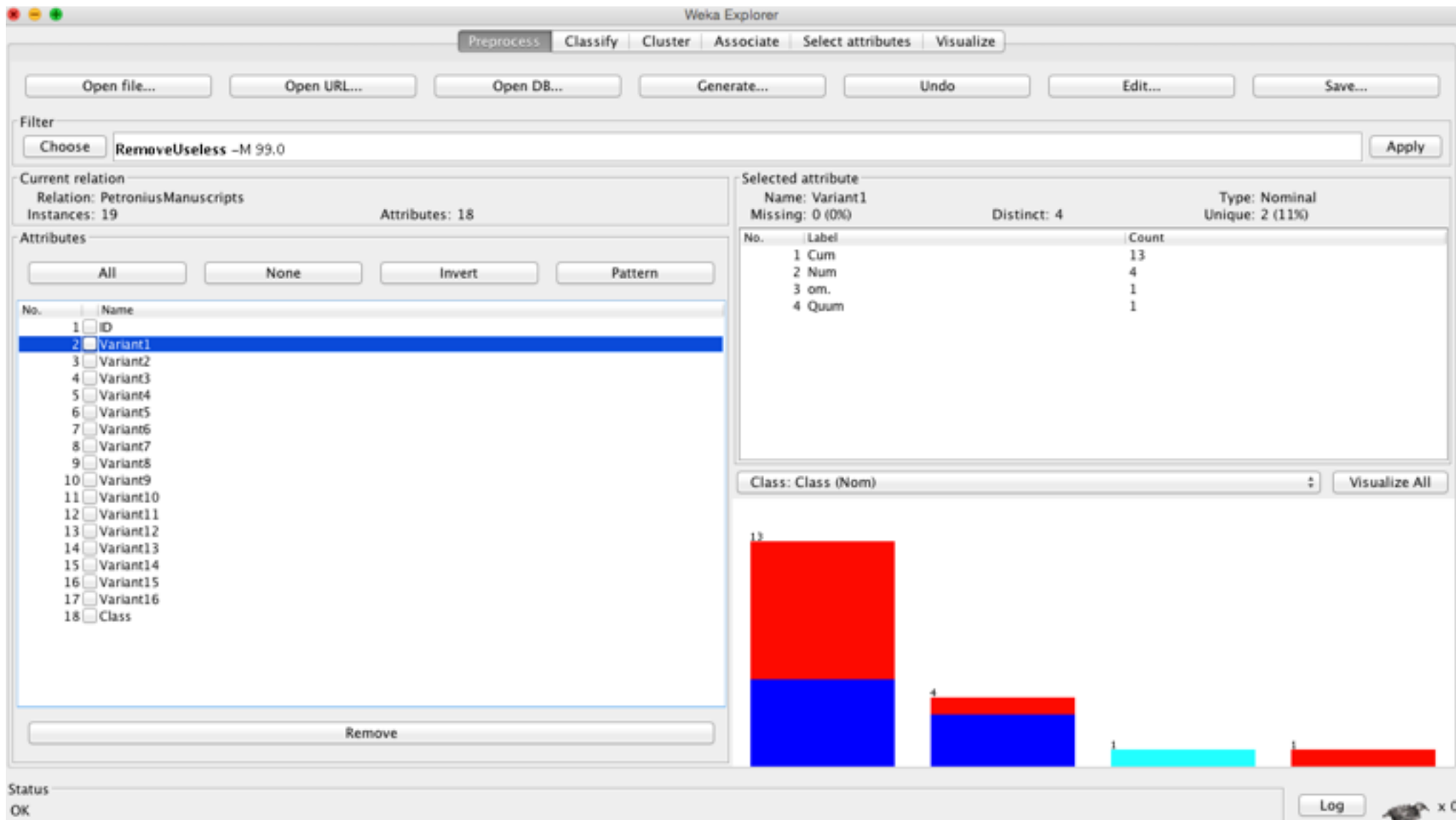
## @data

1,Num,Non,Essent,ituris,facerent,rerum,et,strepitum,orbem\_terrarum,quia,his,et,sed,pestilentiam,omnium,quidem,O

2,Num,Non,Essent,ituris,facerent,rerum,et,strepitum,terrarium\_orbem,quia,iis,et,sed,pestilentiam,omnium,quidem,L

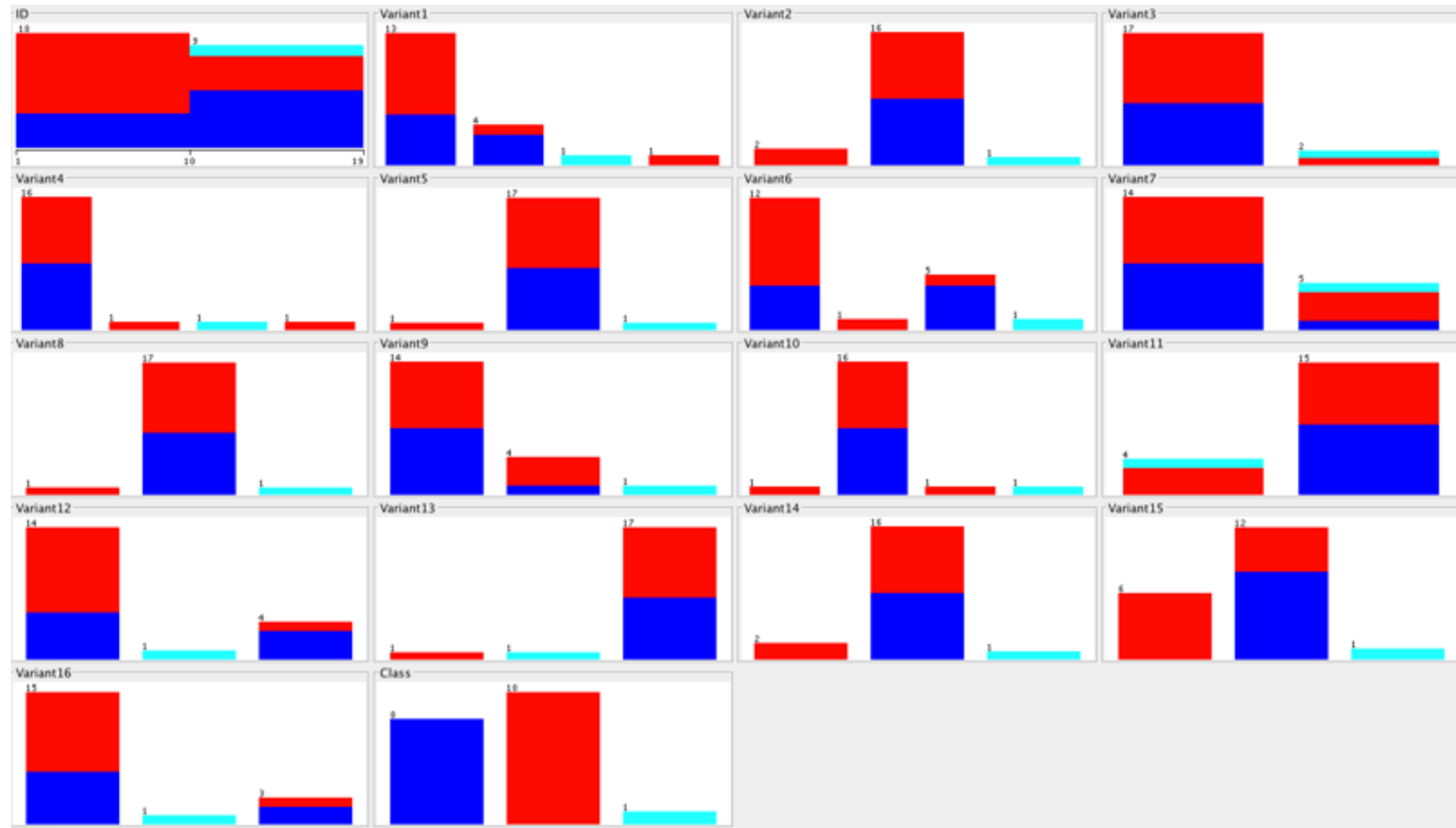
3,Num,Non,Essent,ituris,facerent,rerum,et,strepitum,terrarium\_orbem,quia,iis,sed,sed,pestilentiam,omnium,quidem,L

# Weka





# Weka



# Weka

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose `AttributeSelectedClassifier -E "weka.attributeSelection.WrapperSubsetEval -B weka.classifiers.bayes.NaiveBayes -F 5 -T 0.01 -R 1 --" -S`

Test options:

- ☐ Use training set
- ☐ Supplied test set (Set...)
- ☒ Cross-validation Folds: 10
- ☐ Percentage split %: 66

More options...

(Nom) Class: [dropdown]

Start Stop

Result list (right-click for options):

- 17:45:58 - rules.ZeroR
- 17:46:08 - rules.OneR
- 17:46:21 - bayes.NaiveBayes
- 17:47:02 - lazy.IBk
- 17:47:23 - trees.J48
- 17:47:47 - trees.J48
- 17:49:36 - trees.DecisionStump
- 17:50:36 - trees.RandomTree
- 17:50:42 - trees.BFTree
- 17:50:48 - trees.FT
- 17:50:53 - trees.J48graft
- 17:51:01 - trees.LADTree
- 17:55:00 - trees.J48
- 17:55:27 - trees.RandomForest
- 17:55:40 - bayes.NaiveBayes
- 18:12:22 - trees.UserClassifier
- 18:21:47 - trees.J48
- 18:27:38 - bayes.NaiveBayes
- 18:41:15 - meta.AttributeSelectedClassifier

Classifier output:

Variant16

quidem	7.0	10.0	1.0
om.	1.0	1.0	2.0
solum	3.0	2.0	1.0
[total]	11.0	13.0	4.0

Time taken to build model: 0 seconds

=== Stratified cross-validation ===  
=== Summary ===

Correctly Classified Instances	15	78.9474 %
Incorrectly Classified Instances	4	21.0526 %
Kappa statistic	0.6103	
Mean absolute error	0.2323	
Root mean squared error	0.3788	
Relative absolute error	61.4422 %	
Root relative squared error	87.8515 %	
Total Number of Instances	19	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	1	0.364	0.667	1	0.8	0.795	L
	0.7	0	1	0.7	0.824	0.833	O
	0	0	0	0	0	0.833	Phi
Weighted Avg.	0.789	0.153	0.807	0.789	0.77	0.817	

=== Confusion Matrix ===

a	b	c	<-- classified as
8	0	0	a = L
3	7	0	b = O
1	0	0	c = Phi

Status: OK

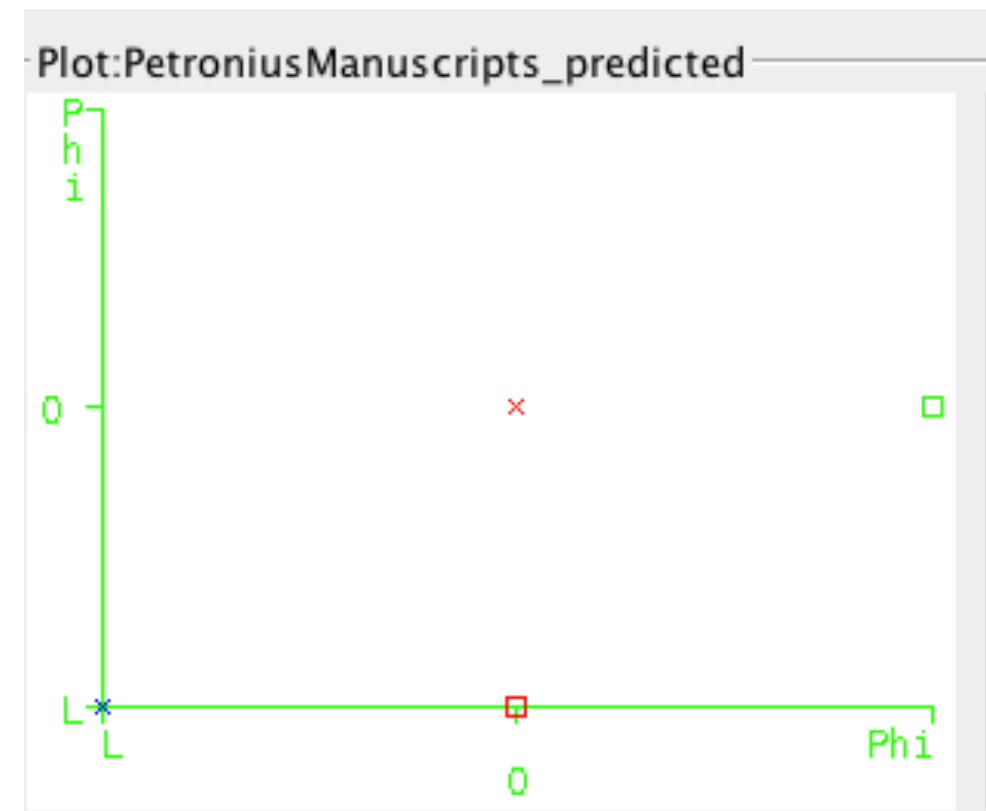
Log x 0

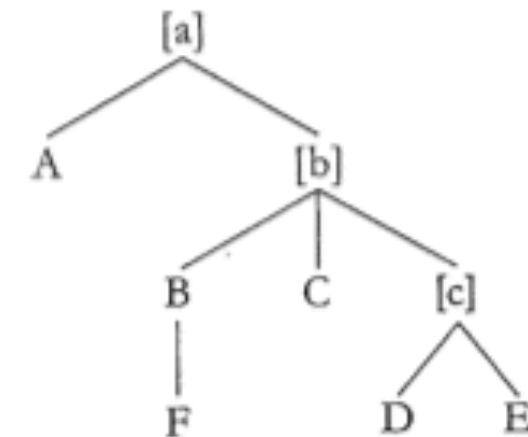
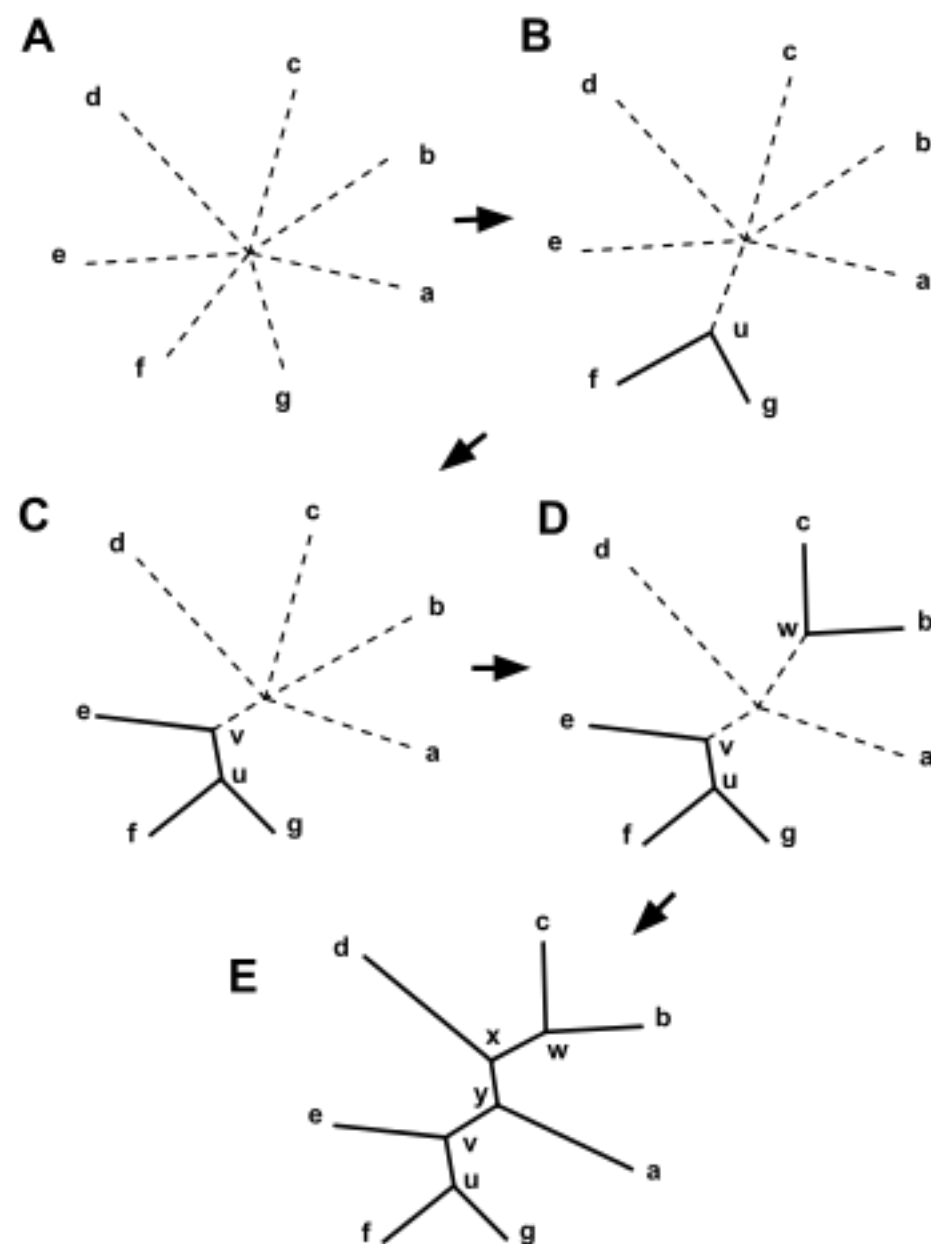


# Weka

$\delta$  = codex a Poggio Florentino anno 1420 repertus, nunc deperditus, ex quo pendent codices novicii qui omnes saeculo XV scripti sunt

- |          |  |
|----------|--|
| $\alpha$ | $\left\{ \begin{array}{l} \mathbf{A} = \text{codicis Parisini lat. 7989 (Orellii } C, \text{ Beck } \textit{Trag. pars prior}) \text{ olim Traguriensis excerpta vulgaria, olim Cipici, cf. } \mathbf{H} \\ \mathbf{I} = \text{codex Indianensis Notre Dame 58} \\ \mathbf{F} = \text{codex Leidensis Vossianus 265 O.81 (Beck } \textit{Vb}) \\ \mathbf{K} = \text{codex Vaticanus lat. 1671} \end{array} \right.$  |
|          | $\left\{ \begin{array}{l} \mathbf{J} = \text{codex Florentinus Laurentianus 37,25 (Beck } \textit{F1}) \\ \mathbf{V} = \text{codex Vindobonensis 179, olim 235 (Endlicheri 218; Beck } \textit{V1}) \\ \mathbf{W} = \text{codex Vindobonensis 3198 (Endlicheri 108), olim Ioannis Sambuci (Beck } \textit{V2, BÜCHELER}^1 \textit{S)} \\ \mathbf{C} = \text{codex Vaticanus Urbinas 670 (Beck } \textit{Vat.)} \\ \mathbf{D} = \text{codex Florentinus Laurentianus 47,31 (Beck } \textit{F2}) \\ \mathbf{G} = \text{codex Guelferbytanus extravag. 299} \\ \mathbf{Q} = \text{codex Vaticanus lat. 3403} \end{array} \right.$ |
|          | $\zeta$  |
|          |  |





Suppose we have species W, X, Y and Z, with the following sequences (in reality, much longer sequences would be used):

W AAAAAAAAA  
X GGAAAAAAAA  
Y CCTTTTAA  
Z CCTTCCAA

The distance matrix would be

	W	X	Y	Z
W	-	2	6	6
X	-	-	6	6
Y	-	-	-	2

and the tree inferred would be:

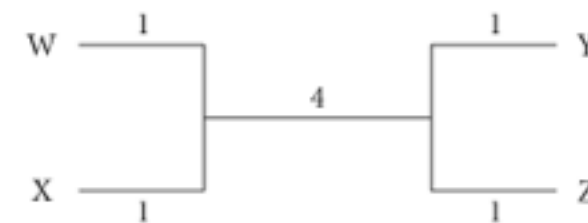


Figure 1. Hypothetical example of a distance matrix analysis.

# Stemmaweb (T. Andrews)

## Stemmaweb - a collection of tools for analysis of collated texts

### Text directory

#### Public text traditions (read-only)

- [Sapientia](#)
- [Parzival 249-255](#)
- [Hanc concordi](#)
- [Quaestiones ad Antiochum ducem \(partially\)](#)
- [Chronicle of Matthew](#)
- [Notre besoin artificiel](#)
- [Parzival artificial](#)
- [Florilegium Coislinianum B](#)

### Text Chronicle of Matthew

- is owned by no one
- is public
- has Armenian as its primary language
- has witnesses: F, J, C, G, Jer, V, A, O, K, D, Y, E, W, L, Z, H, I, B, X

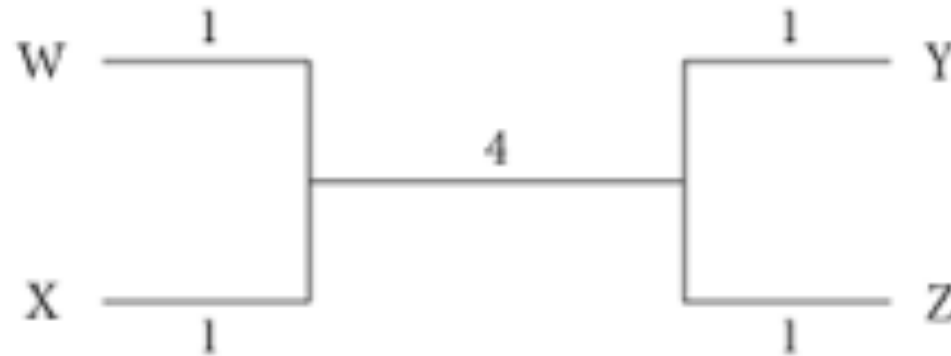
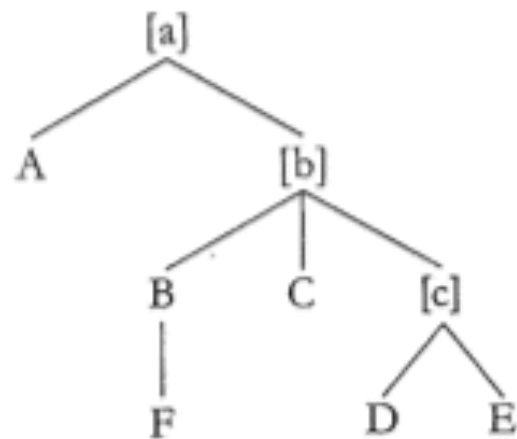


▶ Examine variants against this stemma

▶ View collation and relationships

▶ Download tradition

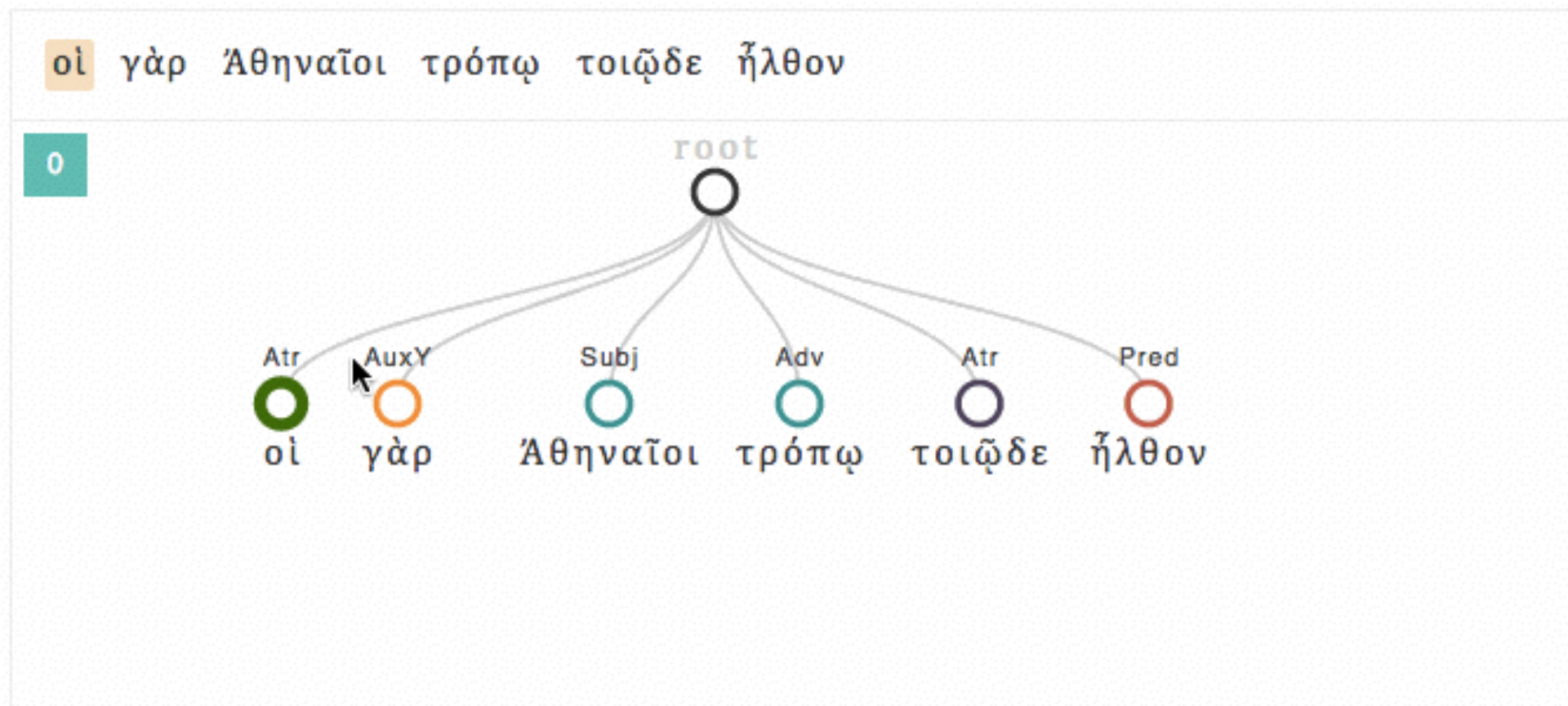
# Scoring Variants



- mathematical “scoring” with the premise that we have full information (can be weighted)
- philological “scoring”: correct/false grammar; meaning; similarity; context -> explanation of the variant, sets the weighting

# Tree Banking & Alignment

# Tree Banking




# Translation Alignment

The screenshot shows a web-based interface for translation alignment. At the top, there is a green tab with the number '0'. Below it, a blue tab with the language code 'grc' is selected. The Greek text is displayed in a serif font: οἱ γὰρ Ἀθηναῖοι τρόπῳ τοιῷδε ἦλθον ἐπὶ τὰ πράγματα ἐν οἷς ηὔξηθησαν .

Below the Greek text, a blue tab with the language code 'en' is selected. The English text is displayed in a sans-serif font: In this way the Athenians came to the circumstances under which they grew in power . The word 'Athenians' is highlighted with an orange background, and a mouse cursor is pointing at it.



# Gradually More Complex

 **ALPHEIOS**

1... 1... 1 >1 >>1 Go to sentence number  Sentence list

Save sentence < Undo Redo > Add Comment Export XML Export Display ☐ Show interlinear text

<p>τὸ μὲν γὰρ ἔτος , ὥς ὠμολογεῖτο , ἐκ πάντων μάλιστα δὴ ἐκείνο ἄνοσον ἐς τὰς ἄλλας ἀσθενείας ἐτύγχανεν ὄν : εἰ δέ τις καὶ προύκαμνέ τι , ἐς τοῦτο πάντα ἀπεκρίθη . τοὺς δὲ ἄλλους ἀπ' οὐδεμῆς προφάσεως , ἀλλ' ἐξαίφνης ὑγιεῖς ὄντας πρῶτον μὲν τῆς κεφαλῆς θέρμαι ἰσχυραὶ καὶ τῶν ὀφθαλμῶν ἐρυθήματα καὶ φλόγωσις ἐλάμβανε , καὶ τὰ ἐντός , ἥ τε φάρυγξ καὶ ἡ γλῶσσα , εὐθὺς αἱματώδη ἦν καὶ πνεῦμα ἄτοπον καὶ δυσῶδες ἠφίει :</p>	<p>Now on the one hand , this year , it is said , out of all the others happened to be especially free of sickness in regards to other diseases : on the other hand , if anyone caught anything , everything turned to this sickness . And others , from no cause , but rather suddenly being healthy before , a strong heat of the head and a redness and burning of the eyes took , and the two inner parts , both the throat and the tongue , immediately began to be bloody and began to throw out a strange and foul breath .</p>
---	--

# Gradually More Complex

Search for documents: Lys. 14 1-47

ήγομαι μέν, ὧ ἄνδρες δικασταί, οὐδεμίαν ὑμᾶς ποθεῖν ἀκοῦσαι πρόφασιν παρὰ τῶν βουλομένων Ἀλκιβιάδου κατηγορεῖν.

selection none 0 unused highlight unused

[ROOT]

PRED AuxK

ήγομαι

AuxY ExD OBJ

μέν ἄνδρες ποθεῖν

AuxX AuxZ ATR AuxX

, ὧ δικασταί, ὑμᾶς ἀκοῦσαι

SBJ OBJ

πρόφασιν παρὰ

ATR ADV

οὐδεμίαν βουλομένων

ATR OBJ

τῶν κατηγορεῖν

morph relation aT search history comments

ἀκοῦσαι 1-11 0 of 18 unused

☒ ἀκούω v--ana--- document  
verb.aor.inf.act

☐ ἀκέω v-pppafn- bsp/morpheus  
verb.pl.pr.part.act.fem.nom

☐ ἀκέω v-pppafv- bsp/morpheus  
verb.pl.pr.part.act.fem.voc

Create new form

# Morphological Parsers

# Morpheus API

```
- <dict>
  <hdwd xml:lang="lat">servo</hdwd>
  <pofs order="1">verb</pofs>
</dict>
- <infl>
  - <term xml:lang="lat">
    <stem>serv</stem>
    <suff>a_</suff>
  </term>
  <pofs order="1">verb</pofs>
  <mood>imperative</mood>
  <num>singular</num>
  <pers>2nd</pers>
  <tense>present</tense>
  <voice>active</voice>
  <stemtype>conj1</stemtype>
  <derivtype>are_vb</derivtype>
</infl>
```

# Topic Modelling



# Bag-Of-Words

- Simplification of texts
- Zellig Harris (1954)
- Word order and/or grammar does not matter



# Topic Modelling

- Collective knowledge continues to grow; it becomes more difficult to find what one is looking for
- Topic Models are more than a search&link-approach
- Zooming in and out is possible with TMs
- Topic Model algorithms do not require prior annotation or labelling of the texts
- Topic Models discover the hidden thematic structure in large archives of documents

# Topic Modelling

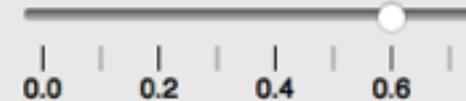
- A method to find clusters of words in large bodies of text
- Those clusters are called topics
- A topic is a recurring pattern of co-occurring words
- Topic models are probabilistic models that are often based on the number of topics in the corpus being assumed and fixed

# Topic Modelling

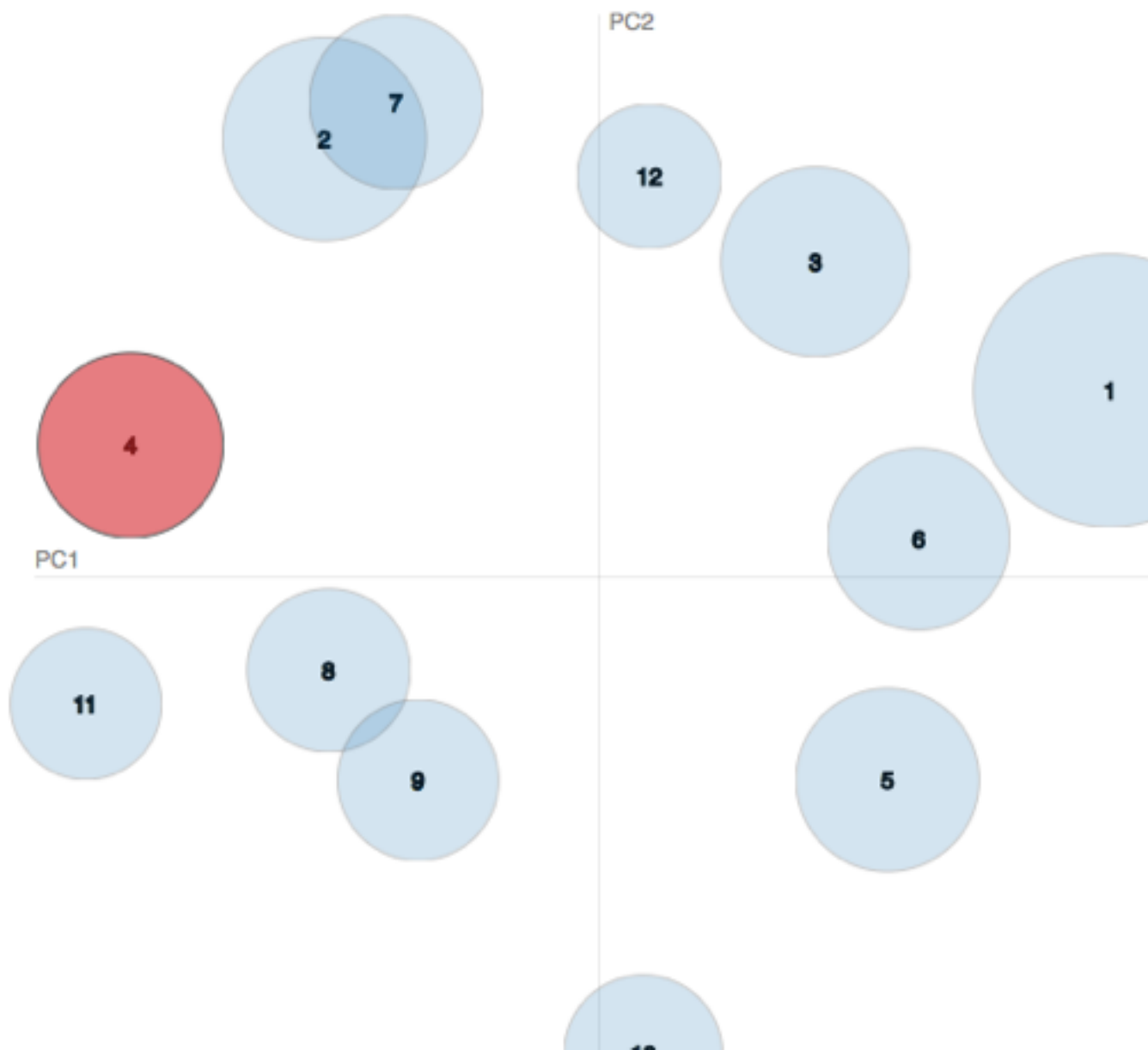
Selected Topic:

Slide to adjust relevance metric:<sup>(2)</sup>

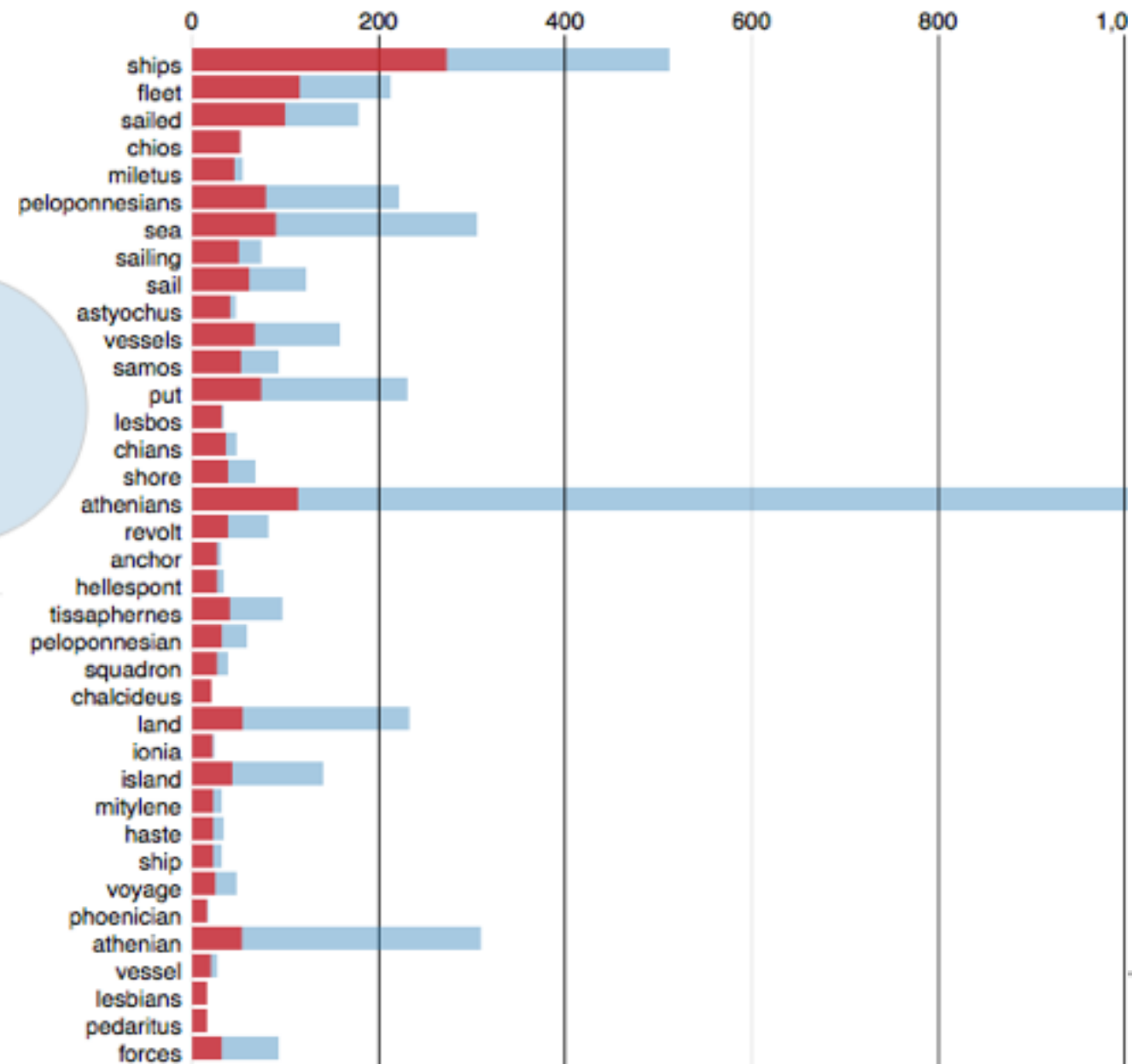
$\lambda = 0.6$



Intertopic Distance Map (via multidimensional scaling)



Top-40 Most Relevant Terms for Topic 4 (8.5% of tokens)



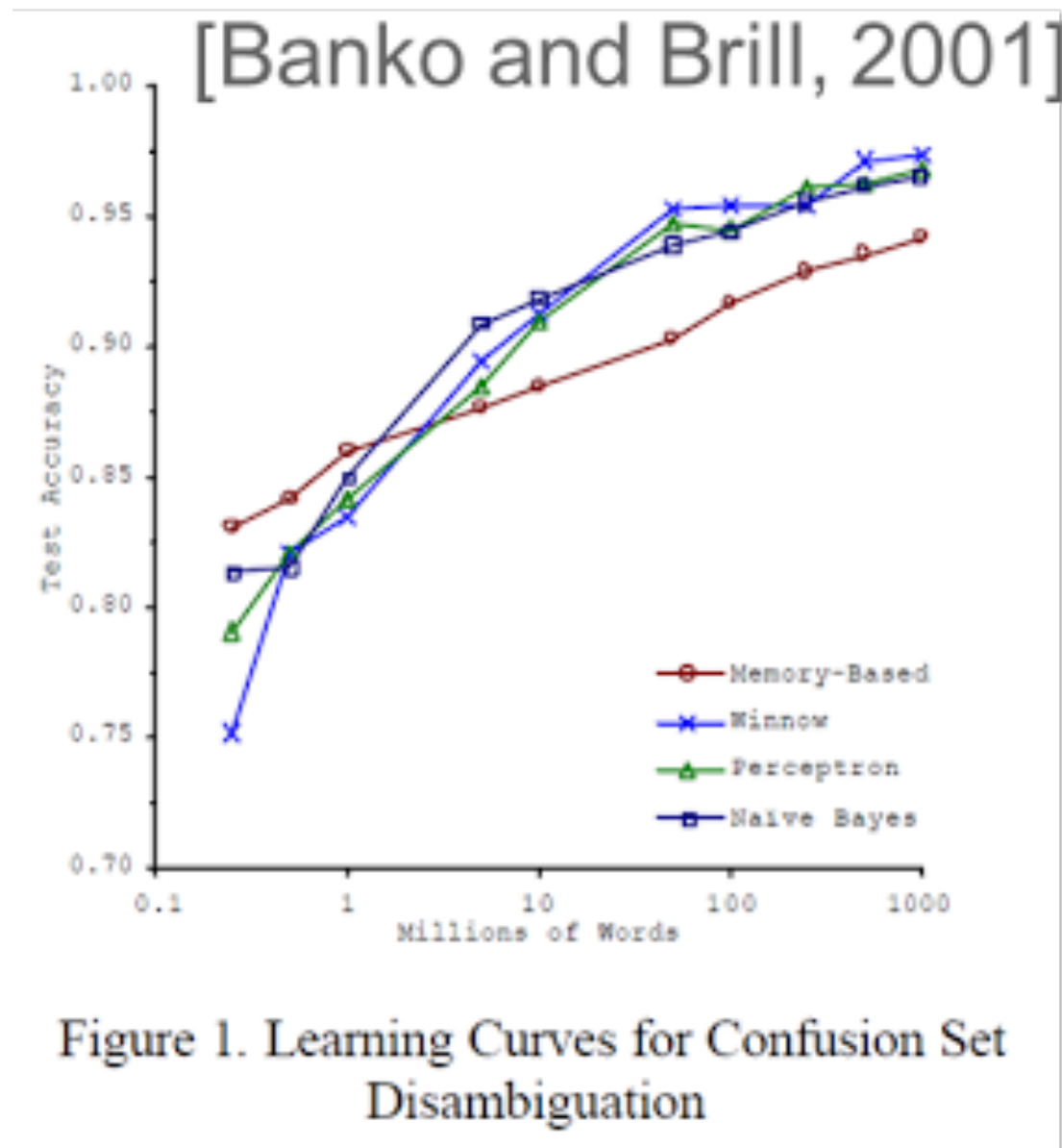


More (better) Data and  
SIMPLE Algorithms  
Often BEAT  
COMPLEX Analytics Models



“We don’t have better algorithms. We just have more data.”

– Peter Norvig (Google)



# Latent-Dirichlet- Allocation (LDA)

# Statistical Inference

Deducing properties of  
an underlying distribution  
by analysis of data.

# LDA

Assumes that one can find a  
Dirichlet Distribution  
by analysing the words in a corpus.

# Dirichlet Distribution

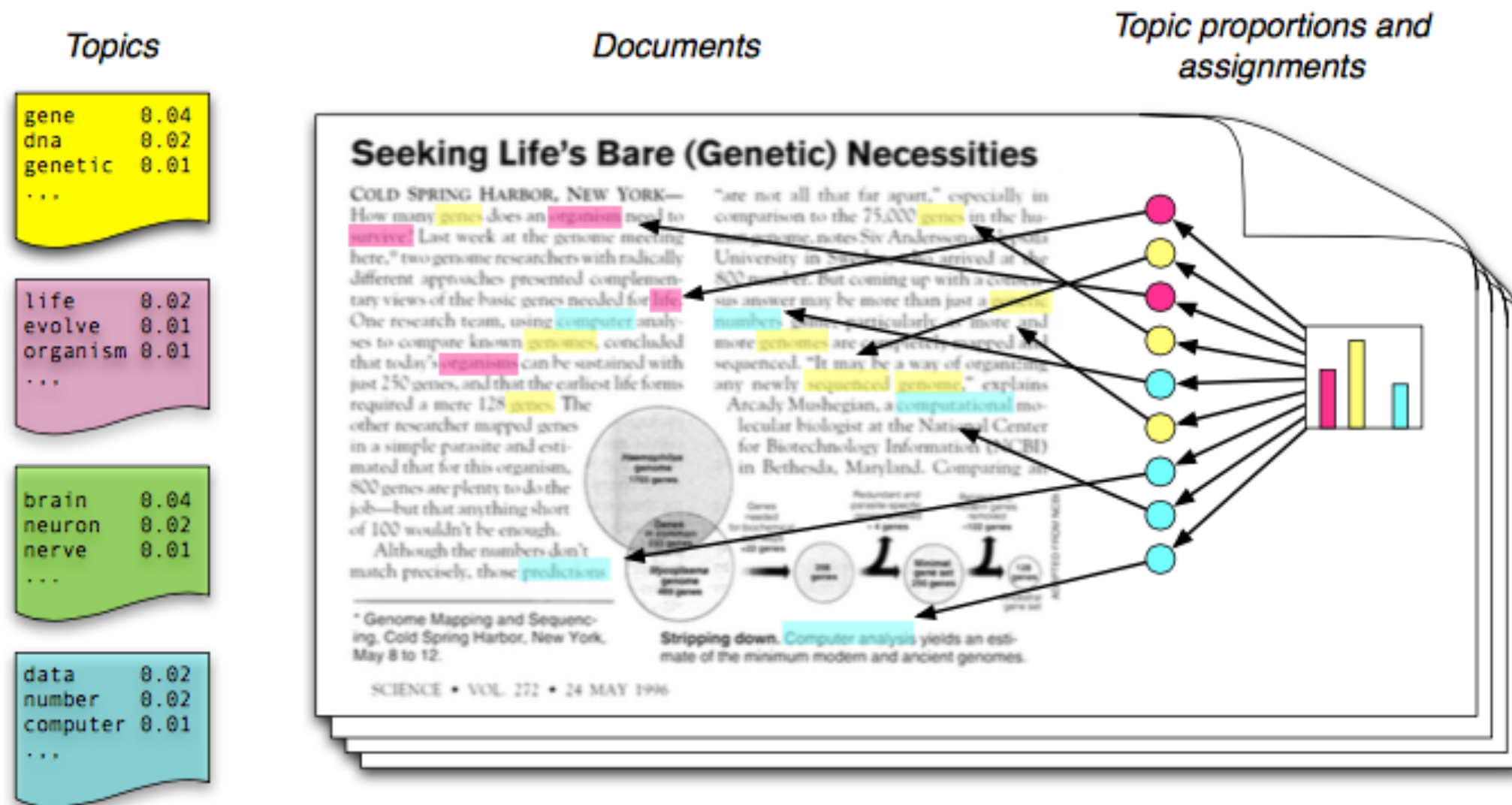
Is a probability distribution  
over all possible  
multinomial distributions



# Dirichlet Distribution of Words in Corpus

1. For each document draw a topic distribution
2. For each word in the document
  - a. Draw a specific topic
  - b. Draw a word

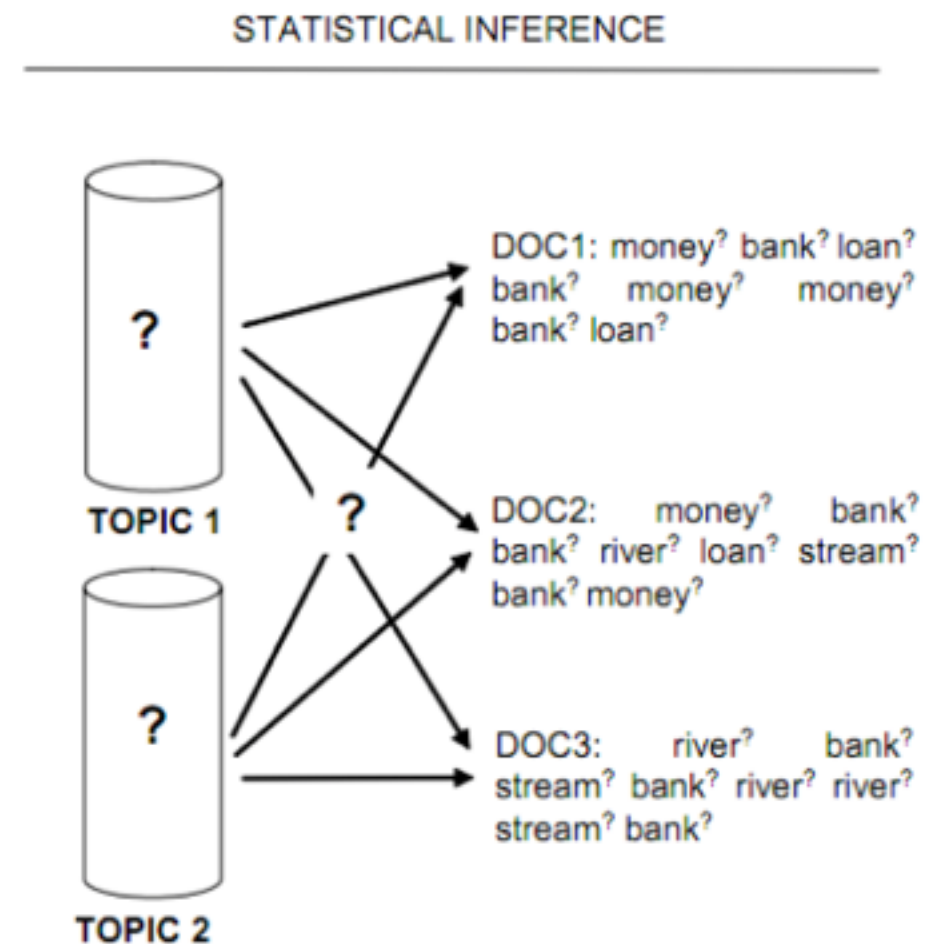
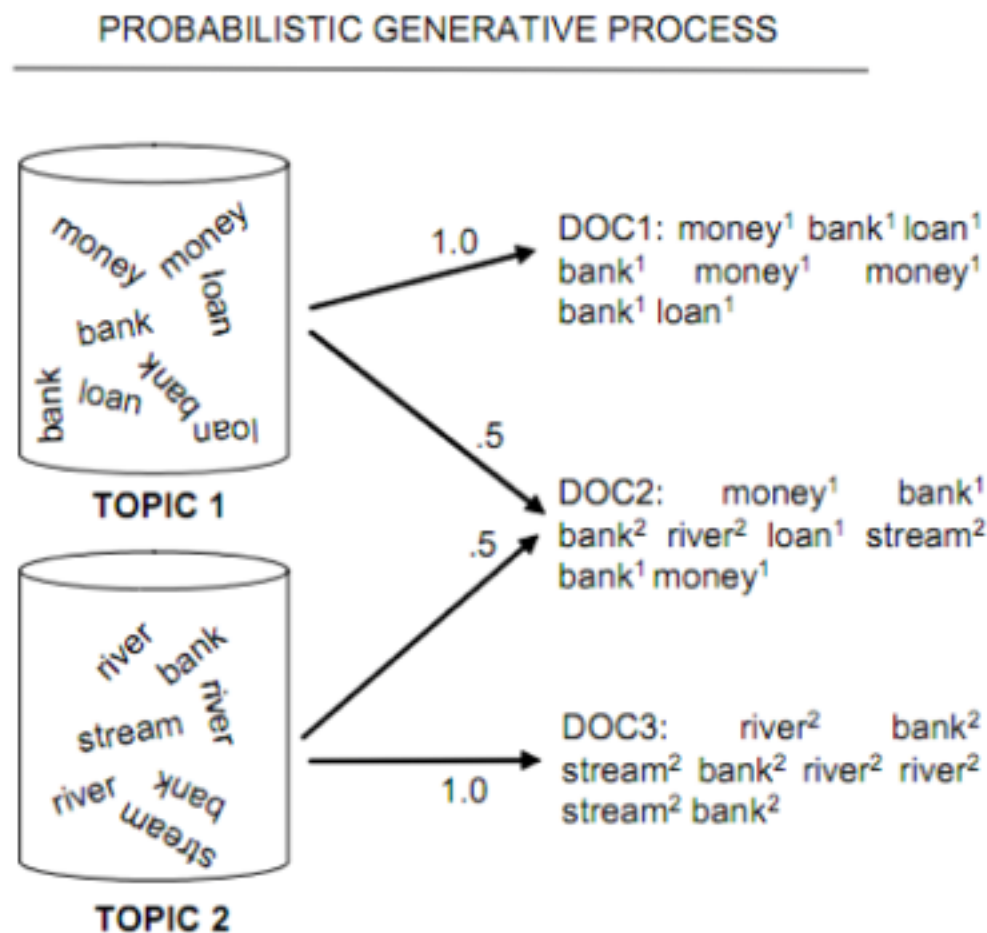
# Dirichlet Distribution of Words in Corpus



# Latent Dirichlet Allocation

- Simplification of how the documents in a dataset were created using Dirichlet Distribution
- Documents are Bag-of-words (order of the words in each document is irrelevant)
- Each corpus has words from a number of known topics
- We do not know which words belong to which topic
- “panel” in document A and “panel” in document B are not the same “word”

# Latent Dirichlet Allocation



# Latent Dirichlet Allocation

- The central goal of topic modelling is to automatically discover the topics from a collection of documents.
- The central inferential problem for LDA is determining the posterior distribution of the latent variables given the document.



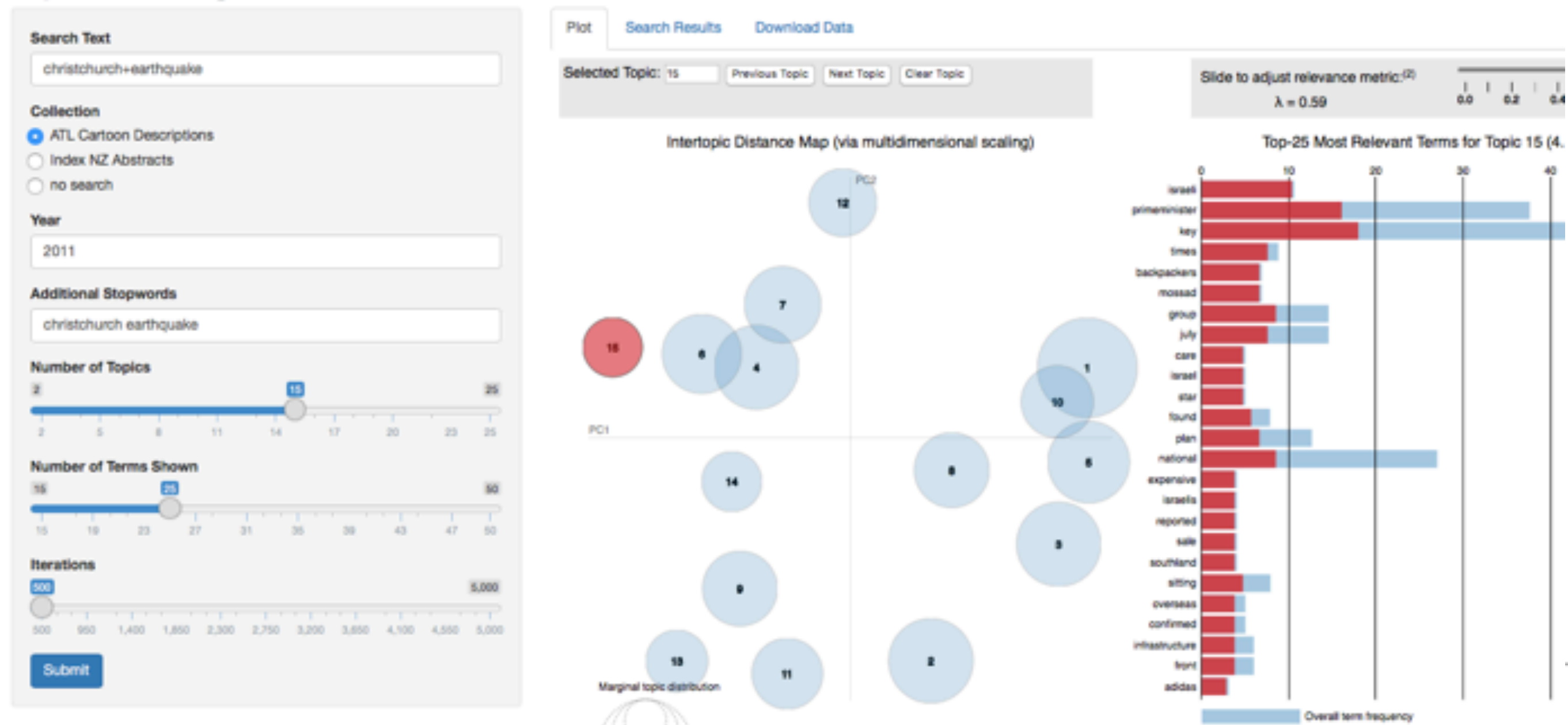
# Success and Results of LDA rely on apriori-set Variables

- Number of topics
- Number of iterations
- Normalisation of the data
- Stopwords

# A Practical Example

- <https://github.com/ThomasK81/TopicModellingR>

## Topic-Modelling based on DNZ API Queries



# What does each variable do?

- Search Text
- Collection
- Year
- Additional Stopwords
- Number of Topics
- Numbers of Terms Shown
- Iterations

# The Risk of Normalisation

- [http://thomask81.github.io/Greek\\_vis/#topic=4&lambda=1&term=](http://thomask81.github.io/Greek_vis/#topic=4&lambda=1&term=)
- ArabicMorph
- ArabicTranslated
- ArabicTM

Topic Modelling is not an end, but a means to an end.



# Sight Reading Finder based on Topic Modelling Thucydides (Tufts, GRK 0103)

- Topic-Modelling Thucydides
- Morphological normalisation of Thucydides  
-> Result: Three different versions of the same text
- Passage-similarity matrix based on TM
- Put it into shiny  
-> Individual App for setting exam questions and enabling preparation.

[https://thomask81.shinyapps.io/sightreading\\_app/](https://thomask81.shinyapps.io/sightreading_app/)

# Questions?

